

# Image Synthesis

EECS 442 – Jeong Joon Park  
Winter 2024, University of Michigan

# Administrivia

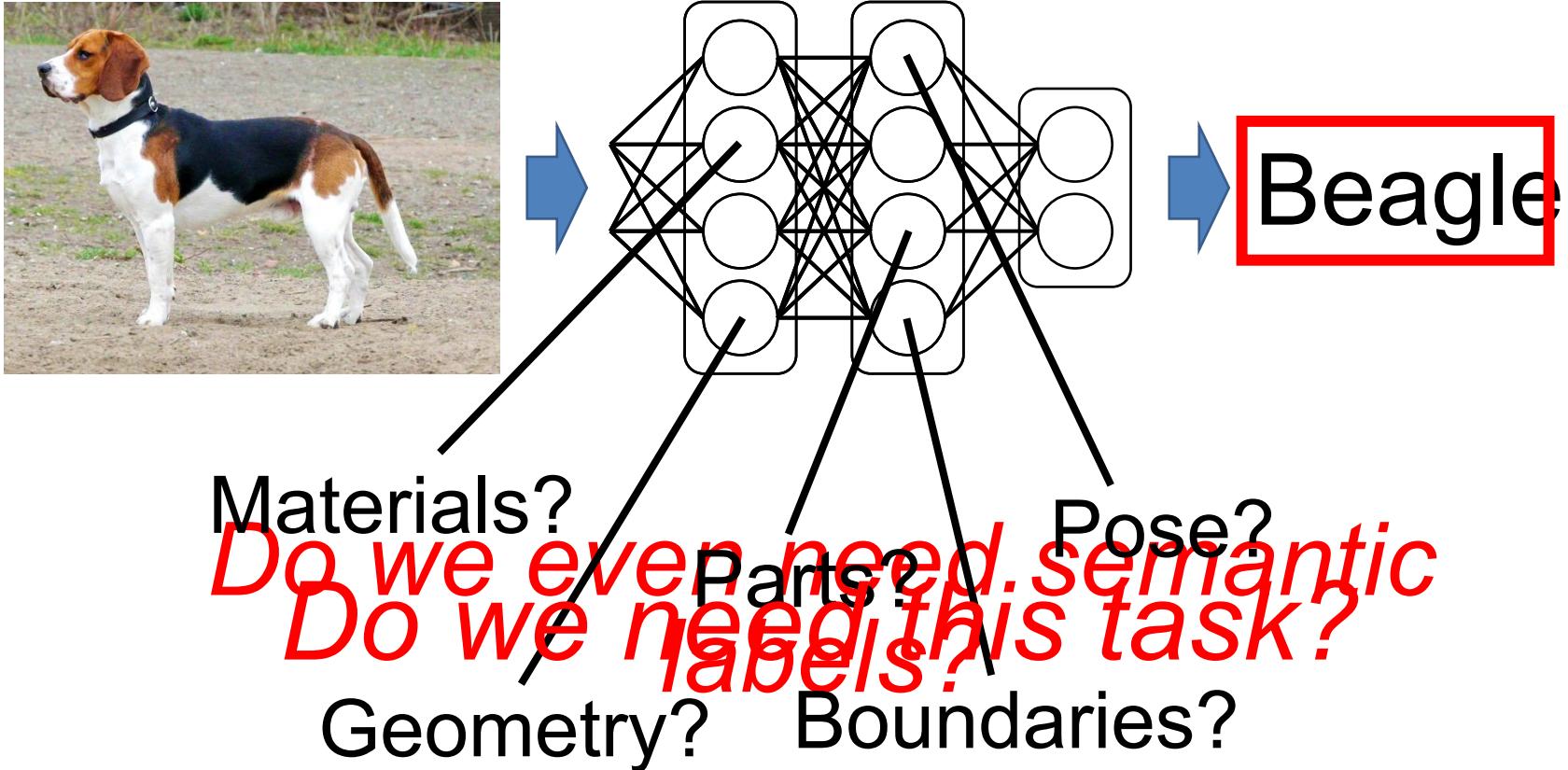
- Mid-term practice exam up – check Piazza
- The real one will cover similar topics
- And easier than the practice one
- Bring a pen – no cheat sheet, laptop, ipad, etc
- Will be seated randomly

# Administrivia

- Project proposal due 27<sup>th</sup>
- We won't give feedback to the proposals
- Instead, come to OH or Discussions to receive feedback

# Recap on Self-Supervised Learning

# ImageNet + Deep Learning



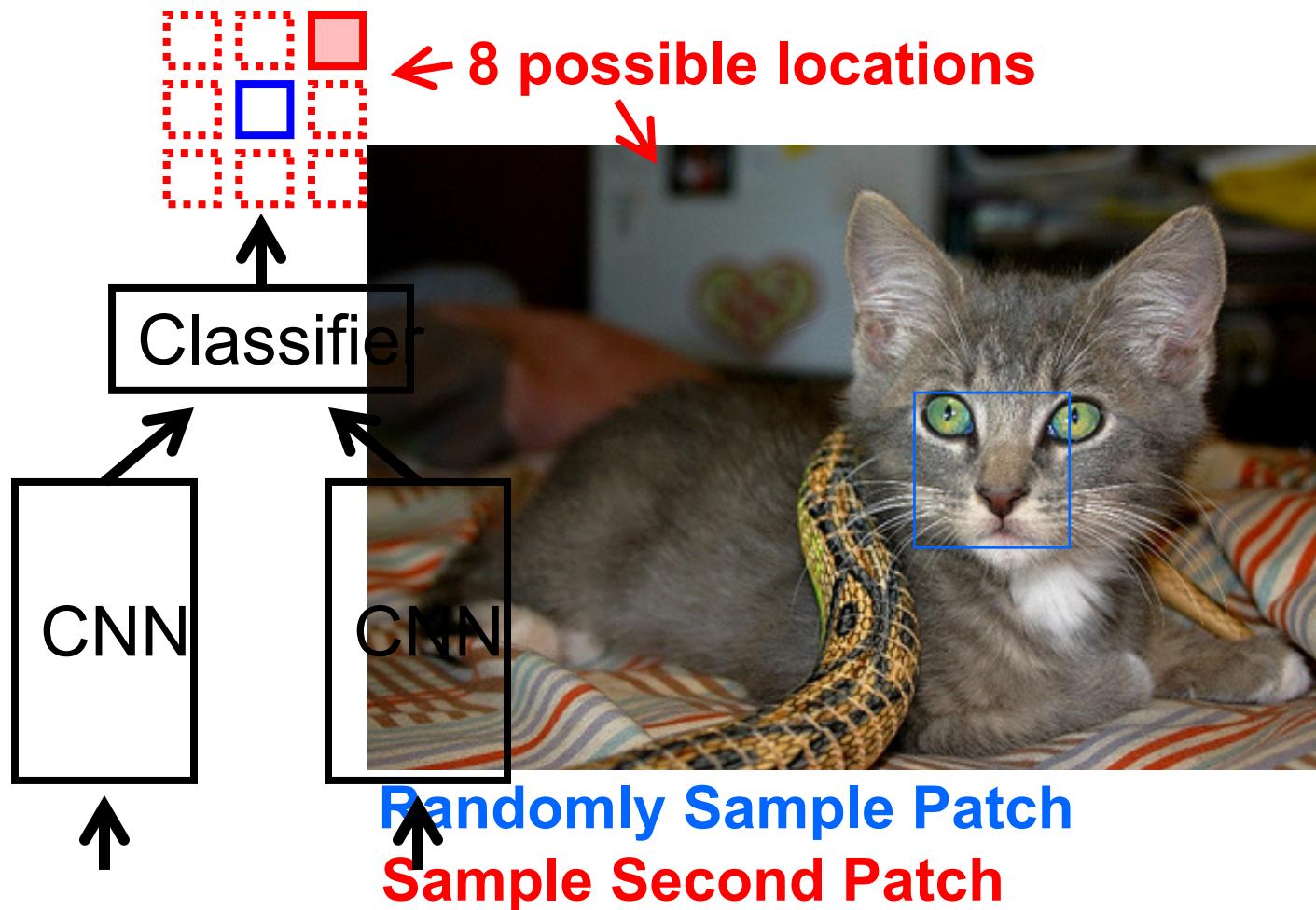
# Context as Supervision

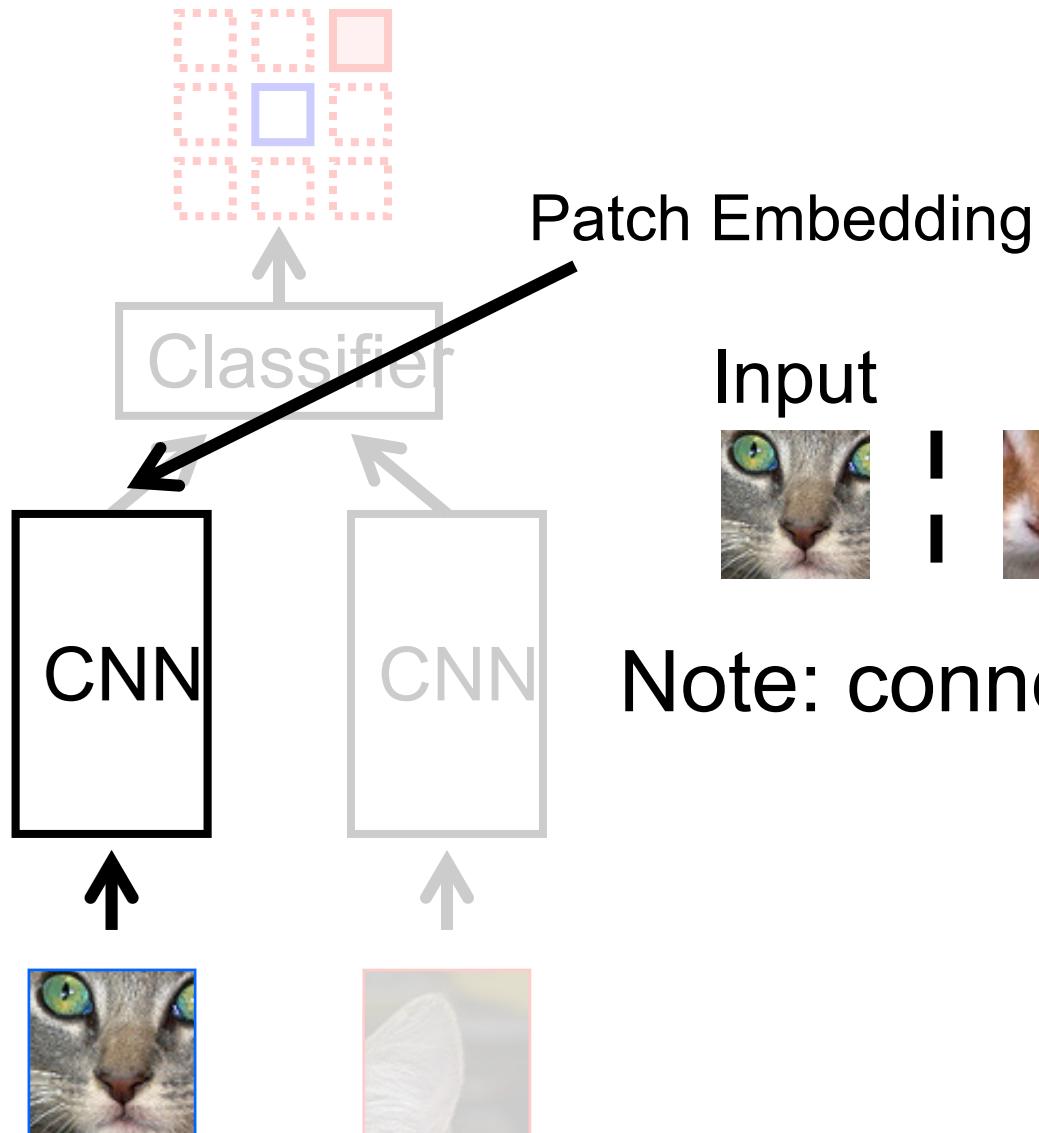
[Collobert & Weston 2008; Mikolov et al. 2013]

house, where the professor lived without his wife and child; or so he said jokingly sometimes: "Here's where I live. My house." His daughter often added, without resentment, for the visitor's information, "It started out to be for me, but it's really his." And she might reach in to bring forth an inch-high table lamp with fluted shade, or a blue dish the size of her little fingernail, marked "Kitty" and half full of eternal milk; but she was sure to replace these, after they had been admired, pretty near exactly where they had been. The little house was very orderly, and just big enough for all it contained, though to some tastes the bric-à-brac in the parlor might seem excessive. The daughter's preference was for the store-bought gimmicks and appliances, the toasters and carpet



# Relative Position Task



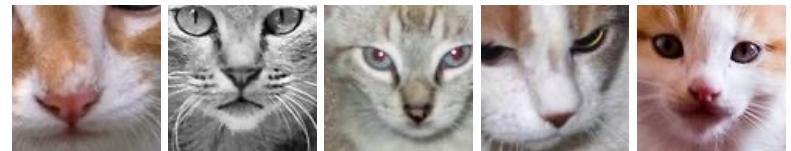


Patch Embedding

Input

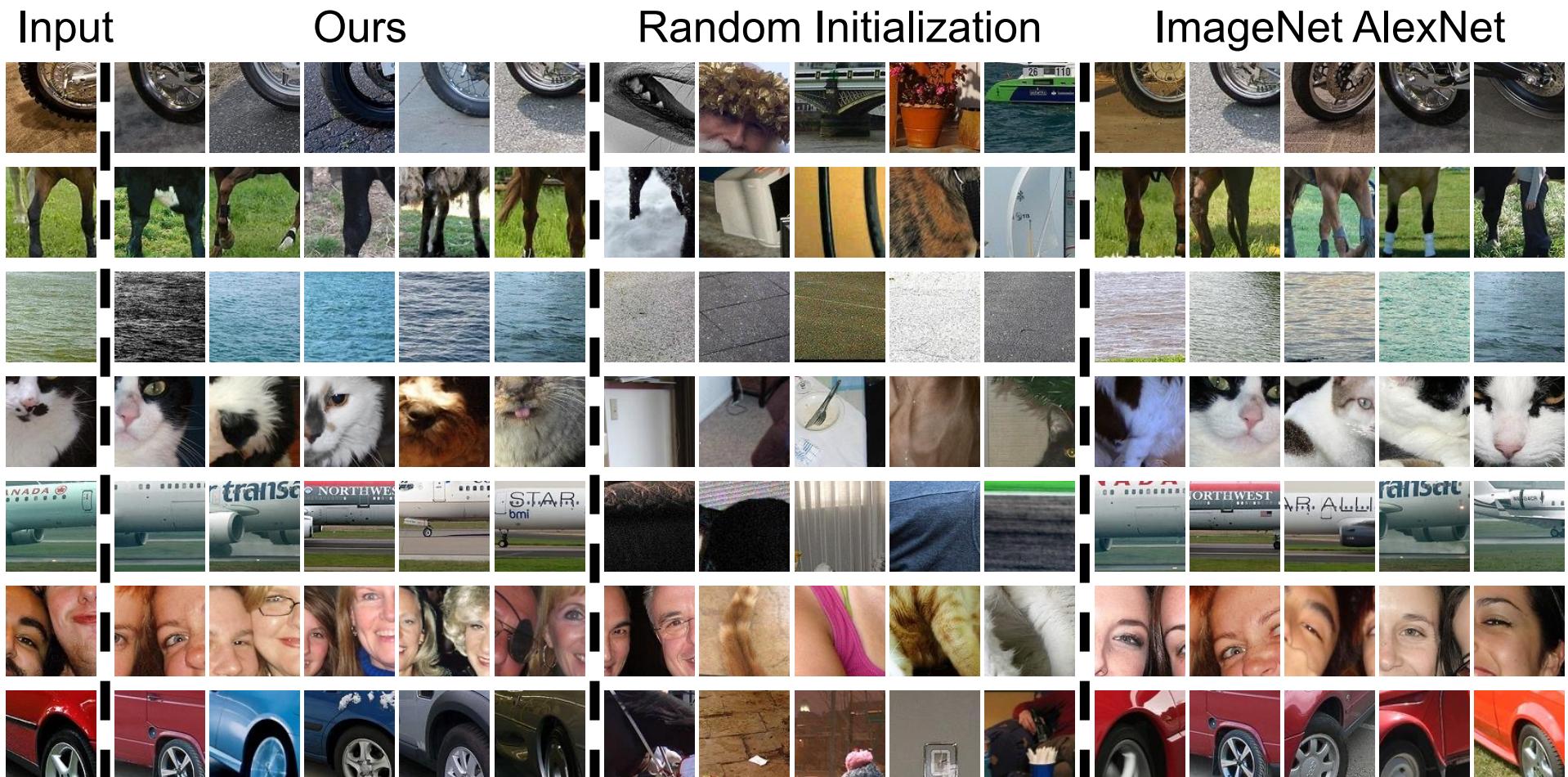


Nearest Neighbors



Note: connects ***across*** instances

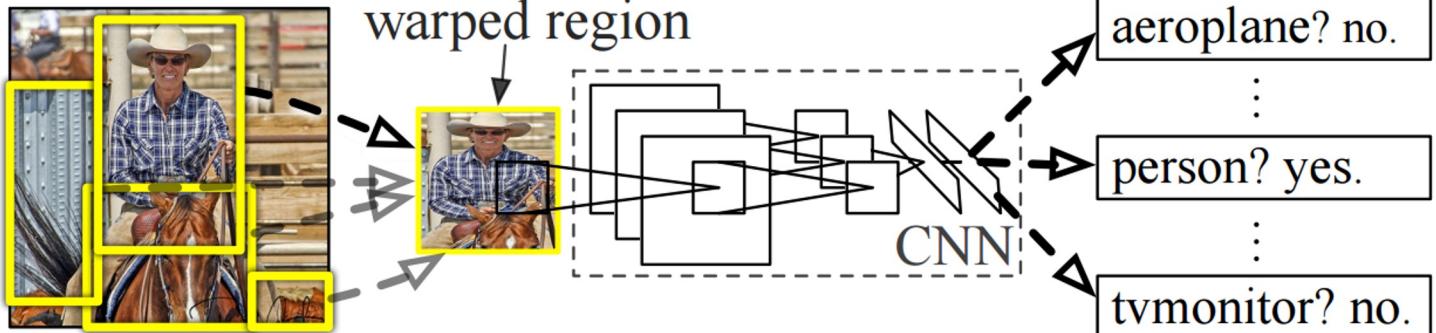
# What is learned?



# Pre-Training for R-CNN



1. Input image



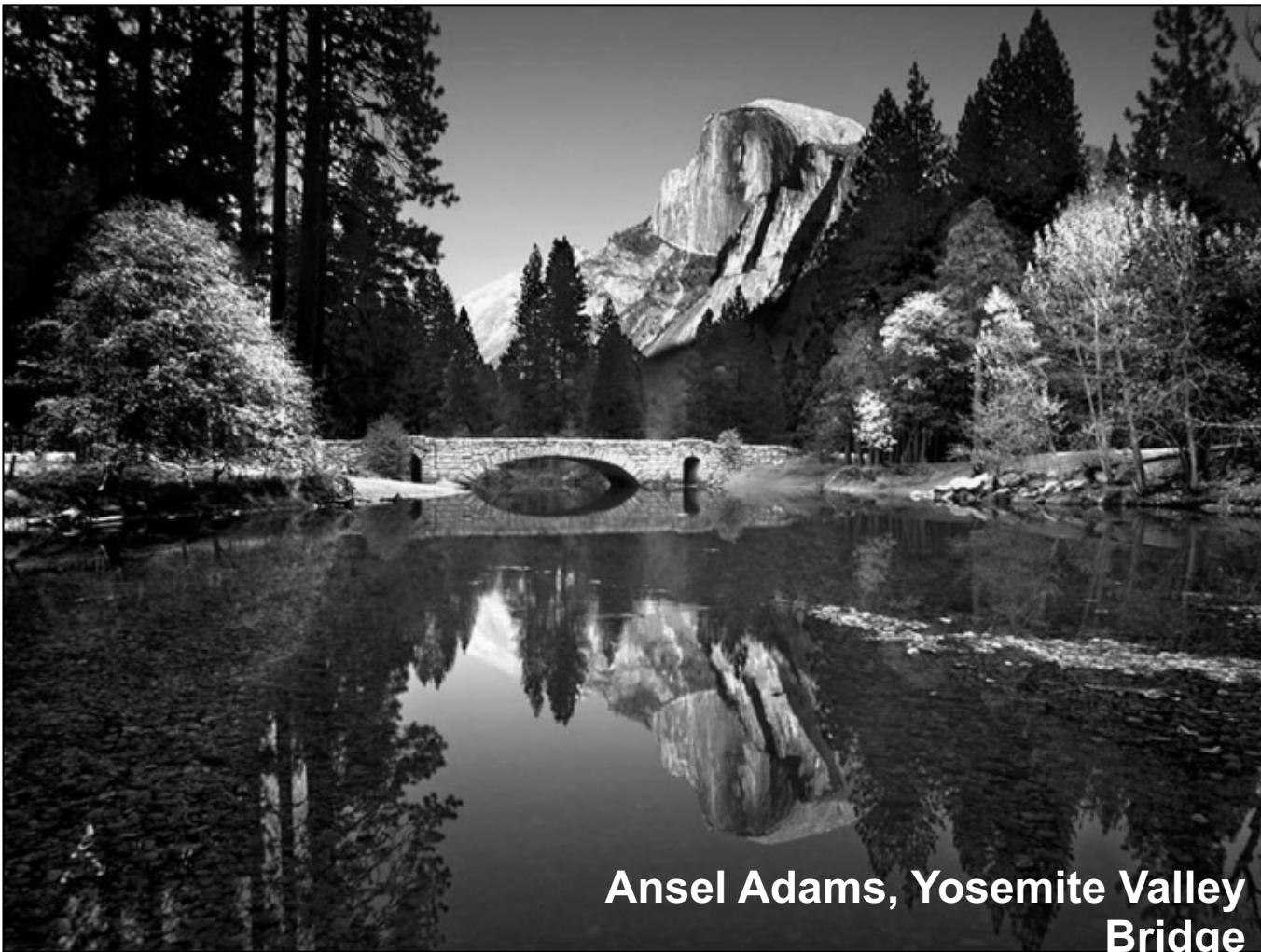
2. Extract region  
proposals (~2k)

3. Compute  
CNN features

4. Classify  
regions

Pre-train on relative-position task, w/o labels

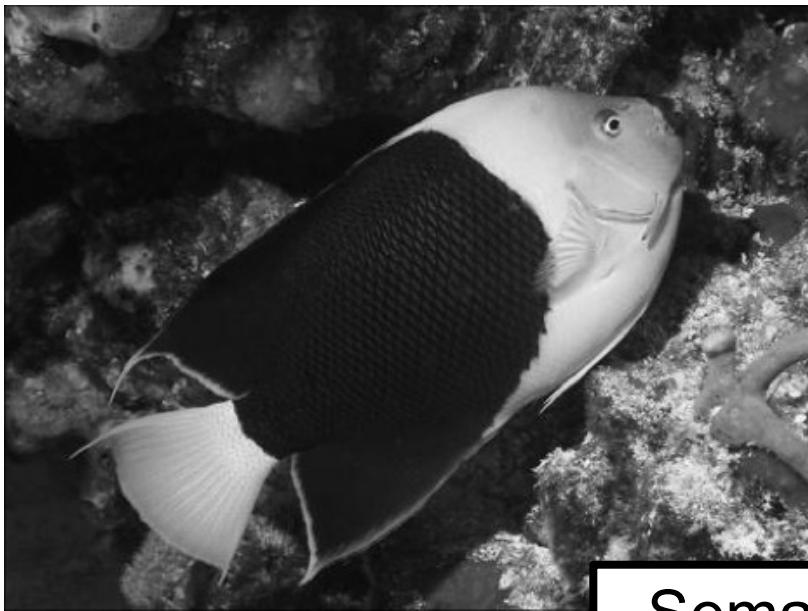
# Other Sources Of Signal



**Ansel Adams, Yosemite Valley  
Bridge**



**Ansel Adams, Yosemite Valley Bridge – Our Result**



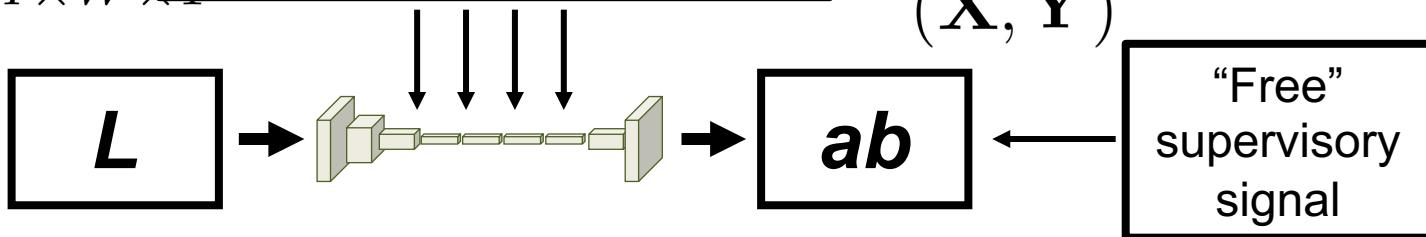
$$\xrightarrow{\mathcal{F}}$$

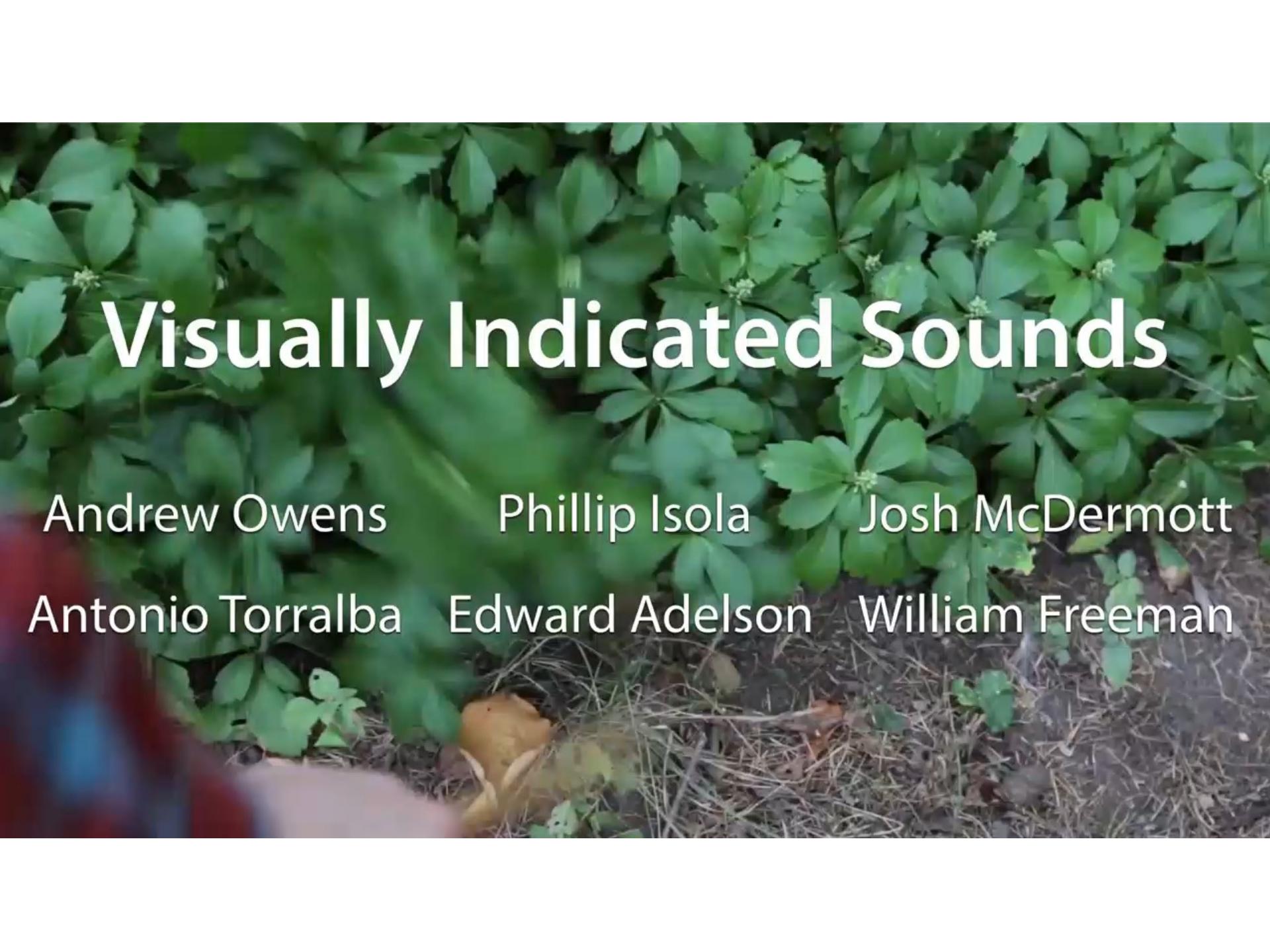


Grayscale image:  $L$   
 $X \in \mathbb{R}^{H \times W \times 1}$

Semantics? Higher-level abstraction?

Concatenate ( $L, ab$ )  
( $X, \hat{Y}$ )





# Visually Indicated Sounds

Andrew Owens

Phillip Isola

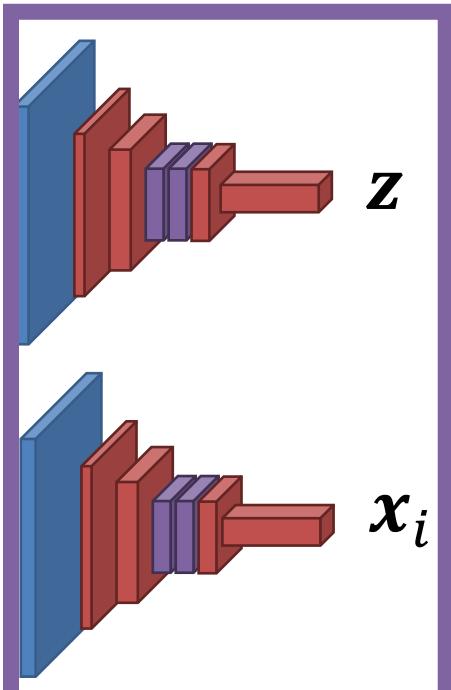
Josh McDermott

Antonio Torralba

Edward Adelson

William Freeman

# Contrastive Learning



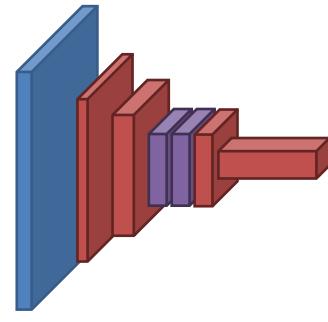
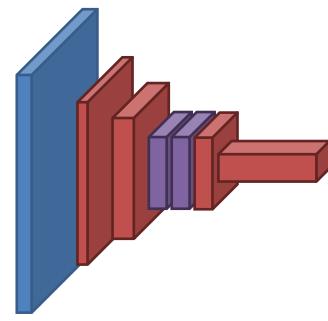
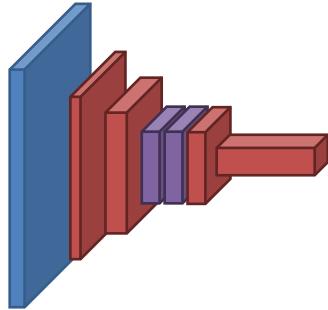
Weights shared  
between networks

Given sample, construct feature  $z$ ; take a bunch of other images, minimize:

$$-\log \left( \frac{\exp(z^T z)}{\exp(z^T z) + \sum_i \exp(z^T x_i)} \right)$$

$$\frac{\exp(z^T z)}{\exp(z^T z) + \sum_i \exp(z^T x_i)} \leq 1$$

# Contrastive Learning



Best performing methods measure distance to augmented sample  $\mathbf{z}'$ :

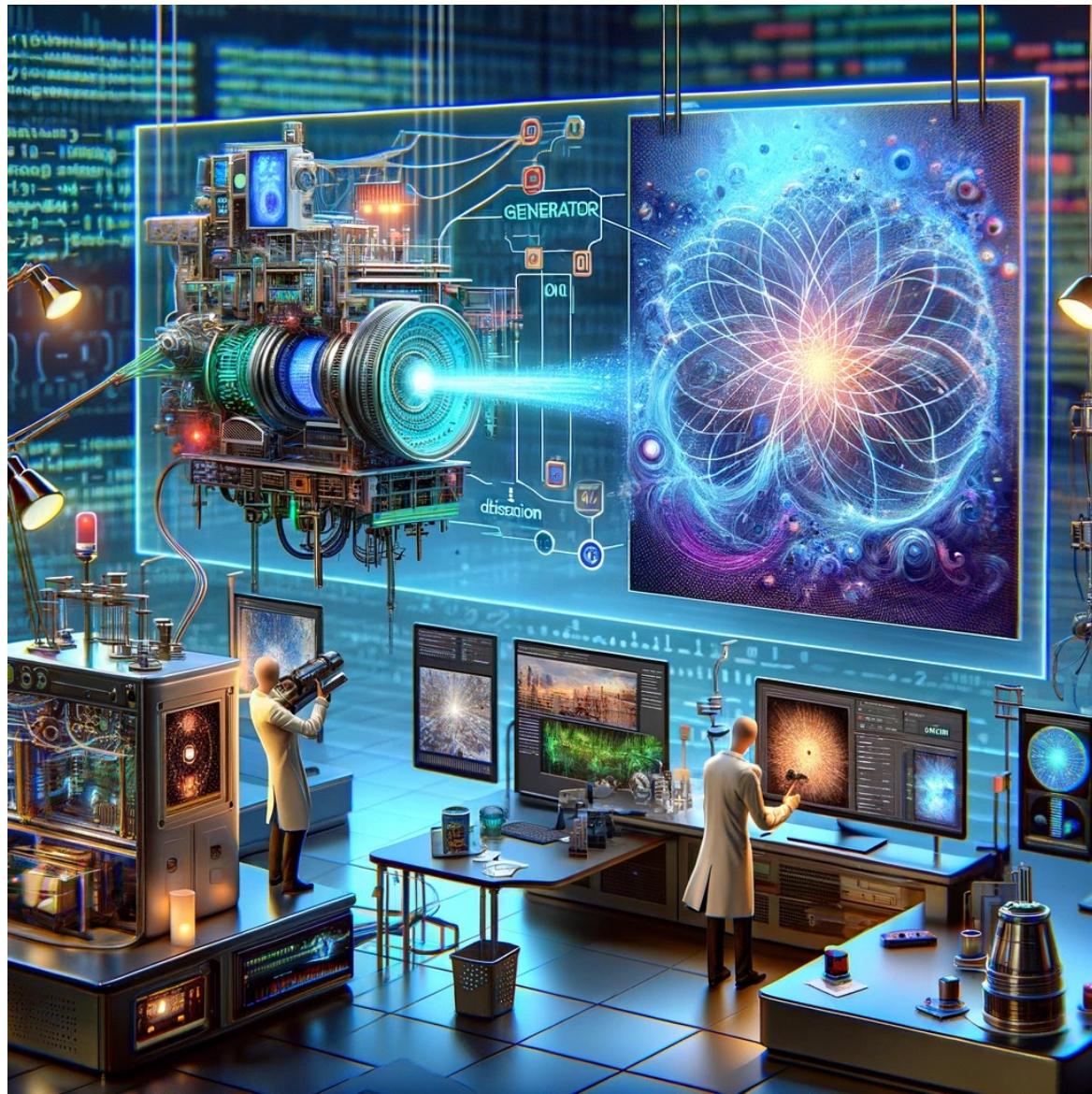
$$-\log \left( \frac{\exp(\mathbf{z}^T \mathbf{z}')}{\exp(\mathbf{z}^T \mathbf{z}') + \sum_i \exp(\mathbf{z}^T \mathbf{x}_i)} \right)$$

Goal: score augmented sample higher than everything else.

# Visual Generative Models

- GANs
- Diffusion Models

→ By ChatGPT



# Discriminative vs Generative Models

**Discriminative Model:**

Learn a probability distribution  $p(y|x)$

**Generative Model:**

Learn a probability distribution  $p(x)$

**Conditional Generative Model:** Learn  $p(x|y)$

**Data: x**



**Label: y**

Cat

# Discriminative vs Generative Models

## Discriminative Model:

Learn a probability distribution  $p(y|x)$

## Generative Model:

Learn a probability distribution  $p(x)$

## Conditional Generative Model:

Learn  $p(x|y)$

Data:  $x$



Label:  $y$

Cat

Probability Recap:

### Density Function

$p(x)$  assigns a positive number to each possible  $x$ ; higher numbers mean  $x$  is more likely

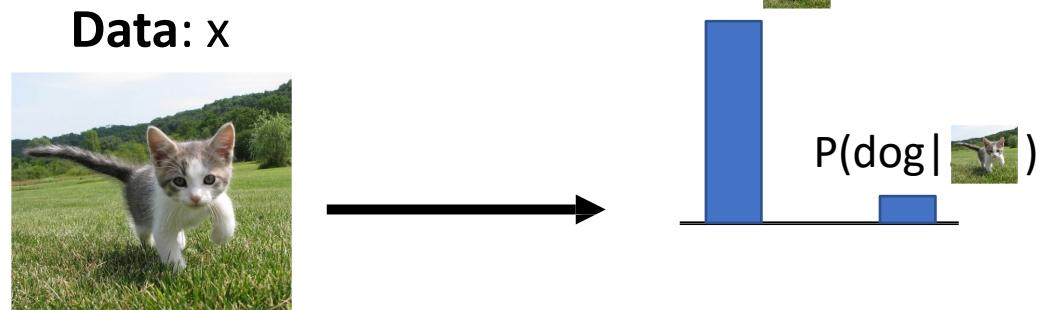
Density functions are **normalized**:

$$\int_X p(x)dx = 1$$

Different values of  $x$  **compete** for density

# Discriminative vs Generative Models

**Discriminative Model:**  
Learn a probability distribution  $p(y|x)$

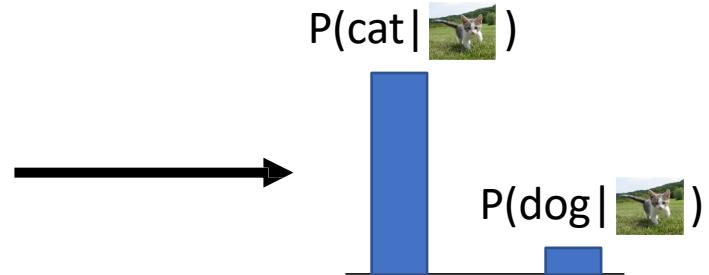


**Generative Model:**  
Learn a probability distribution  $p(x)$

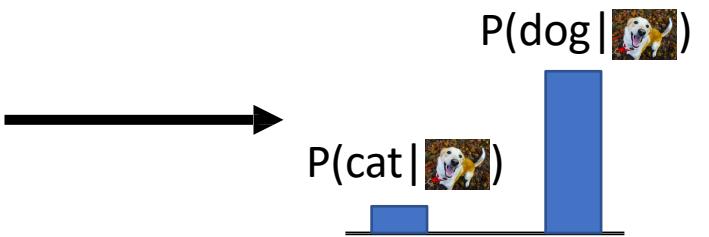
**Conditional Generative Model:** Learn  $p(x|y)$

# Discriminative vs Generative Models

**Discriminative Model:**  
Learn a probability distribution  $p(y|x)$



**Generative Model:**  
Learn a probability distribution  $p(x)$

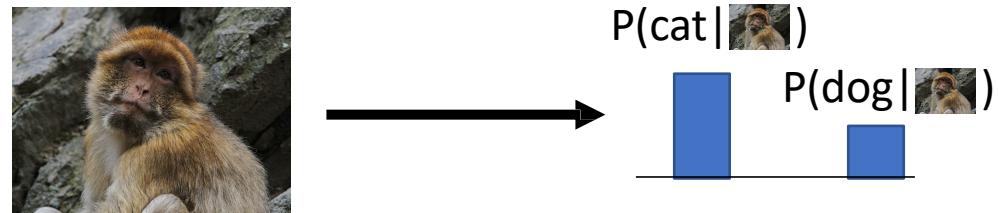


**Conditional Generative Model:** Learn  $p(x|y)$

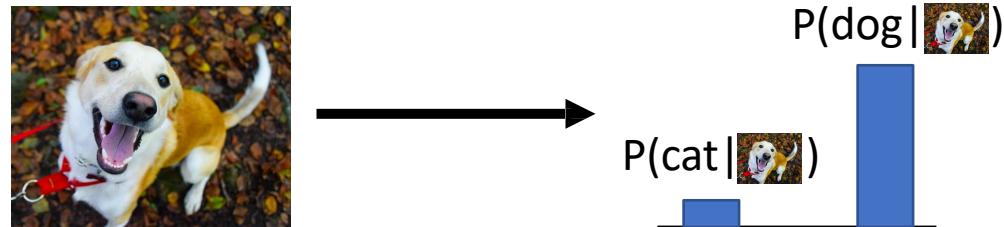
Discriminative model: the possible labels for each input "compete" for probability mass.  
But no competition between images

# Discriminative vs Generative Models

**Discriminative Model:**  
Learn a probability distribution  $p(y|x)$



**Generative Model:**  
Learn a probability distribution  $p(x)$

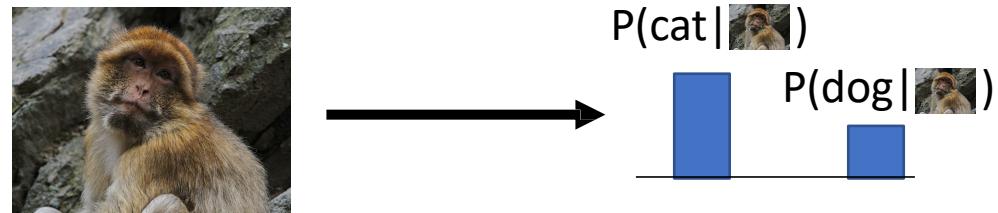


**Conditional Generative Model:** Learn  $p(x|y)$

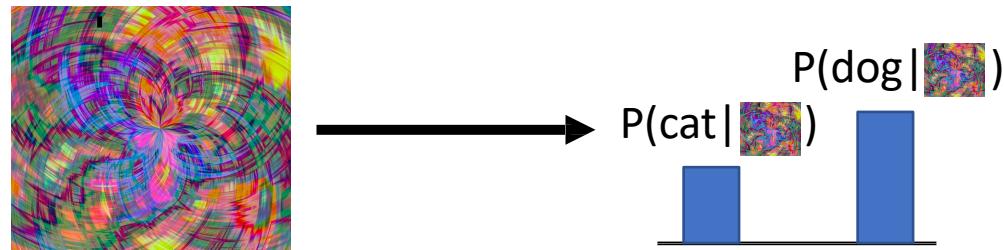
Discriminative model: No way for the model to handle unreasonable inputs; it must give label distributions for all images

# Discriminative vs Generative Models

**Discriminative Model:**  
Learn a probability distribution  $p(y|x)$



**Generative Model:**  
Learn a probability distribution  $p(x)$



**Conditional Generative Model:** Learn  $p(x|y)$

Discriminative model: No way for the model to handle unreasonable inputs; it must give label distributions for all images

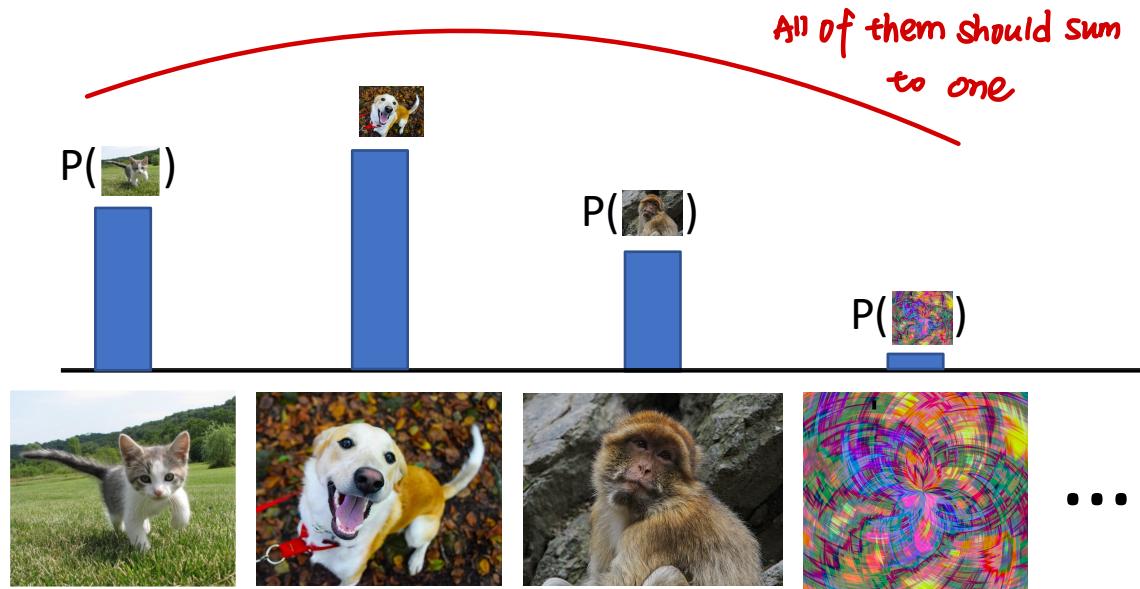
# Discriminative vs Generative Models

## Discriminative Model:

Learn a probability distribution  $p(y|x)$

**Generative Model:**  
Learn a probability distribution  $p(x)$

**Conditional Generative Model:** Learn  $p(x|y)$



Generative model: All possible images compete *(unlikely)*. with each other for probability mass

Requires deep image understanding! Is a dog more likely to sit or stand? How about 3-legged dog vs 3-armed monkey?

Model can “reject” unreasonable inputs by assigning them small values

# Discriminative vs Generative Models

**Discriminative Model:**

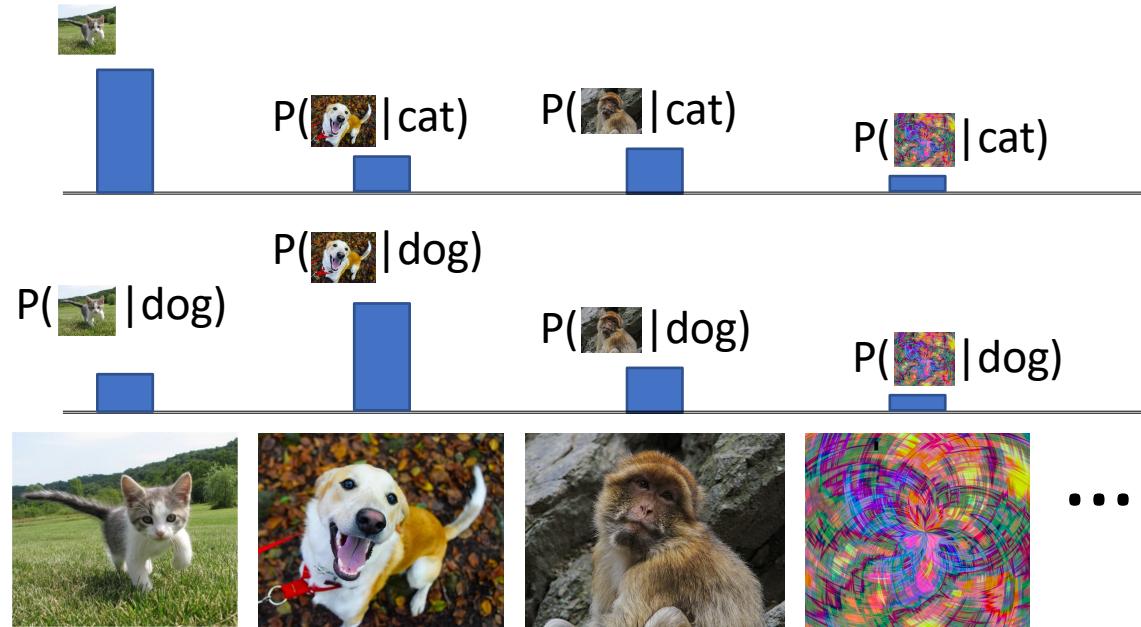
Learn a probability distribution  $p(y|x)$

**Generative Model:**

Learn a probability distribution  $p(x)$

**Conditional Generative Model:** Learn  $p(x|y)$

↑  
given the label.



Conditional Generative Model: Each possible label induces a competition among all images

# Discriminative vs Generative Models

**Discriminative Model:**

Learn a probability distribution  $p(y|x)$

**Generative Model:**

Learn a probability distribution  $p(x)$

Recall Bayes' Rule:

$$P(x | y) = \frac{P(y | x)}{P(y)} P(x)$$

**Conditional Generative Model:** Learn  $p(x|y)$

# Discriminative vs Generative Models

**Discriminative Model:**

Learn a probability distribution  $p(y|x)$

**Generative Model:**

Learn a probability distribution  $p(x)$

**Conditional Generative Model:** Learn  $p(x|y)$

Recall Bayes' Rule:

$$P(x | y) = \frac{P(y | x)}{P(y)} P(x)$$

Conditional Generative Model

Discriminative Model      (Unconditional) Generative Model

Prior over labels

*your intuition or some statistic of the*

We can build a conditional generative *world say* model from other components, etc. *that there are more dogs' labels than rats' labels.*

*(prior distribution of the world)*

# Progress in Generative Models of Images



2014



2015



2016



2017



2018

Slide credit: Ian Goodfellow, 2019



StyleGAN2

<https://github.com/NVlabs/stylegan3>



StyleGAN3 (Ours)

[Karras et al., “Alias-Free Generative Adversarial Networks”, 2021]

# Images and Text

TEXT PROMPT

an armchair in the shape of an avocado....

AI-GENERATED  
IMAGES



Edit prompt or view more images↓

$P(\text{image} \mid \text{caption})$

TEXT PROMPT

a store front that has the word 'openai' written on it....

AI-GENERATED  
IMAGES



# Inverse Problems with GenAI

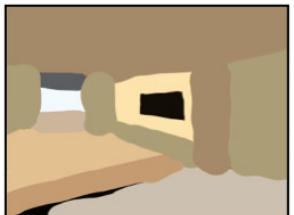
Stroke Painting to Image



Input

Output

Stroke-based Editing

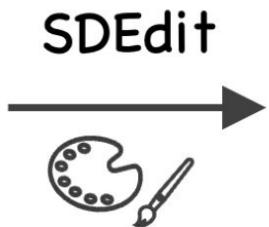


# Inverse Problems with GenAI

Which image is real?



User  
 @StefanoErmon



Output

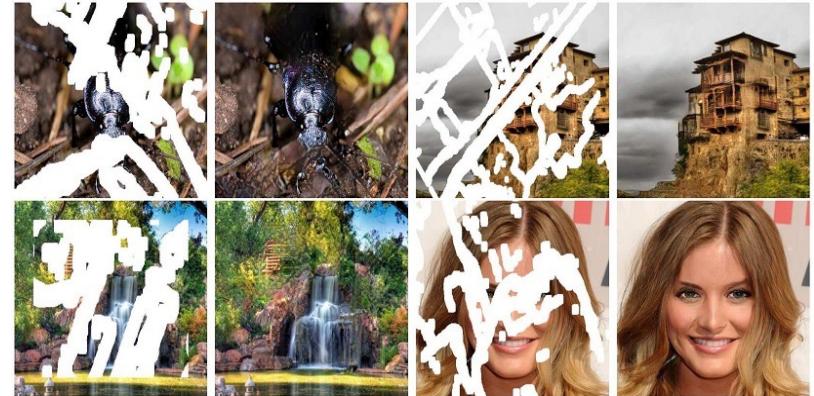
# Inverse Problems with GenAI

P(high resolution | low resolution)



Menon et al, 2020

P(full image| mask)



Liu al, 2018

P(color image| greyscale)



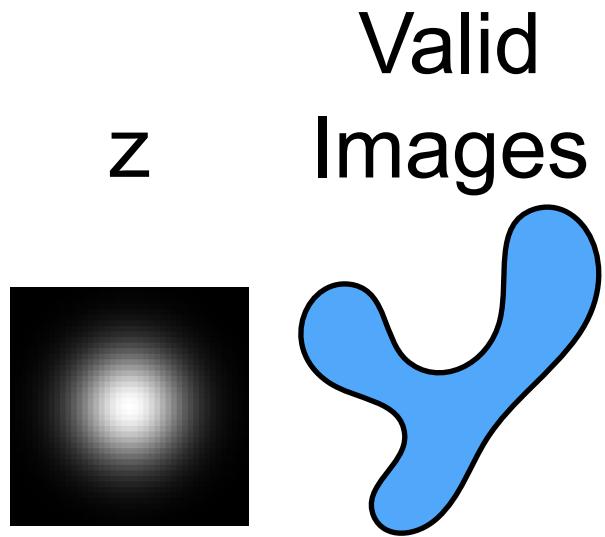
Antic, 2020

# How Many Images Are There?

- Set height and width to 1024
- Assume  $256^3$  (aka  $2^{24}$ ) values per pixel
- **How many images can I create?**
- $(256^3) \sim 1M$
- **Why might it be quite a bit less?**

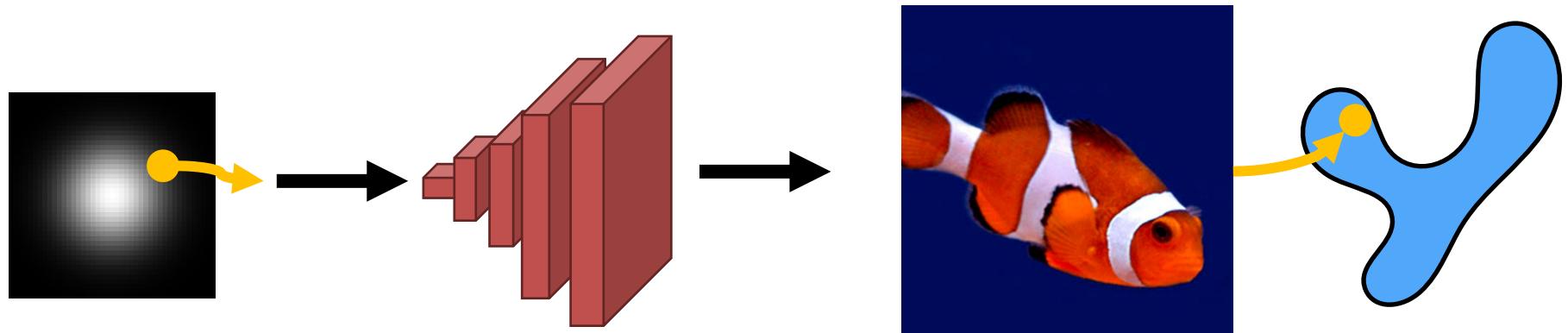
# Learning a Mapping

- Want to learn a mapping
- **From:** a “latent” space  $z$ , often assume to be the result of sampling N-D Gaussian noise
- **To:** the space of valid images



# Generating Data

$z$        $G(\cdot)$        $x$       Valid Images



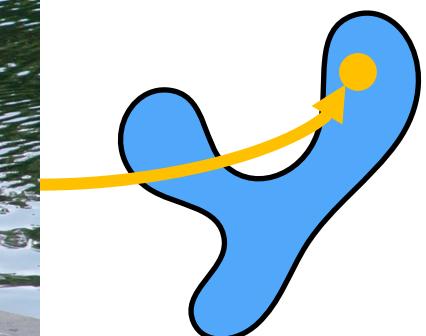
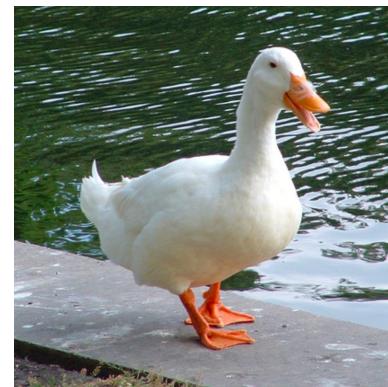
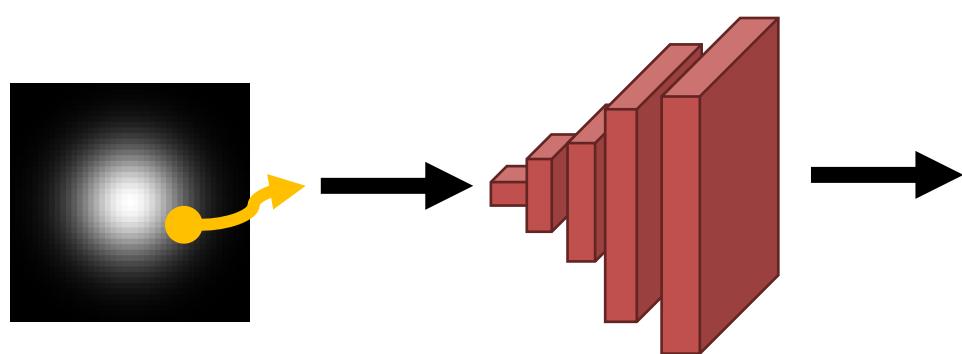
# Generating Data

$z$

$G(\cdot)$

$x$

Valid  
Images

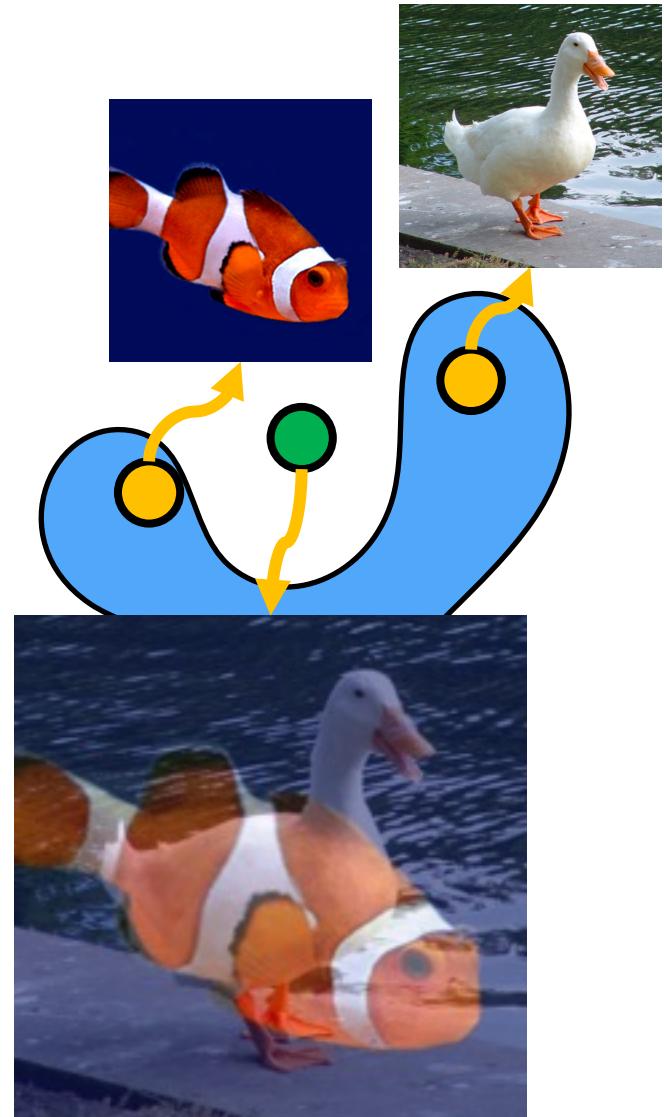


# Why The Funny Shape?

Given two valid images, what about their average?

Key things to remember

- Linear combinations of images aren't valid images. (*off-manifold*)
- Explains funny shape (“manifold”) and need for a deep network

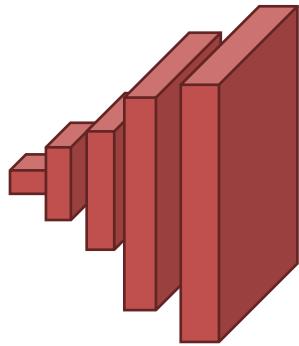


# Generative Adversarial Networks

- Generator tries to make fake images – accepts noise and makes an image
- Discriminator tries to identify fakes – outputs  $p(\text{fake}^{\text{real}})$

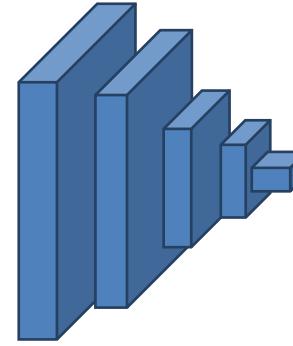
Generator

$z \quad G(\cdot)$



*feed both generated image & real image*

Discriminator  
 $D(\cdot)$



Real vs  
fake

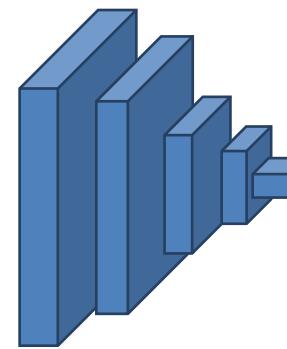
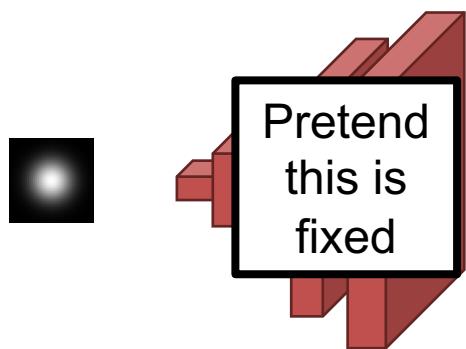
# Generative Adversarial Networks

$z$

$G(\cdot)$

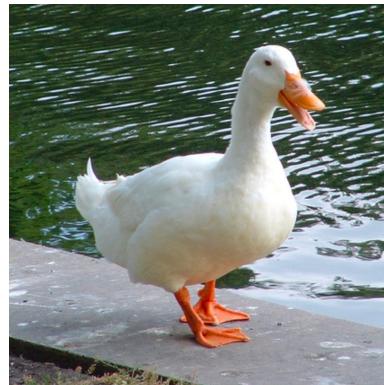
$G(z)$

$D(\cdot)$

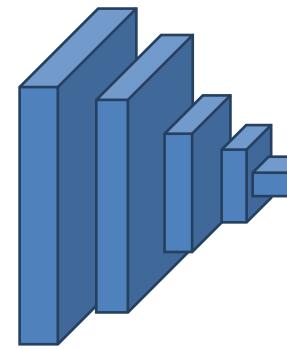


Want: fake  
(low)

$x$



$D(\cdot)$



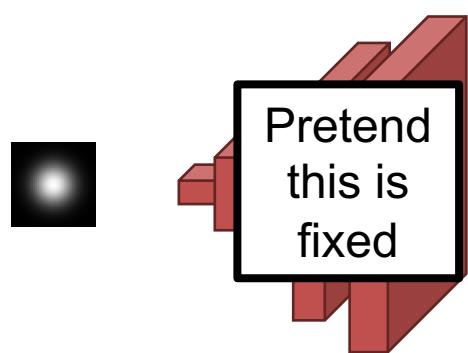
Want: real  
(high)

$D(x) =$   
 $p(x \text{ is real})$   
Want to maximize  
 $\log p(\text{real})$  for real  
data

$$\arg \max_D E_{z,x} [\log (1 - D(G(z))) + \boxed{\log(D(x))}]$$

# Generative Adversarial Networks

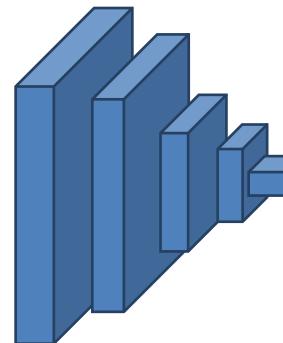
$z$        $G(\cdot)$



$G(z)$



$D(\cdot)$



Want: fake  
(low)

$Gz$  is real

$$1 - D(G(z)) = \underbrace{p(G(z))}_{\text{is fake}}$$

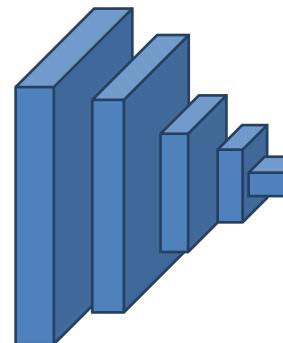
Want to maximize  
 $\log p(\text{fake})$  for  
fake data

*Because we want  
the model to learn.*

$x$



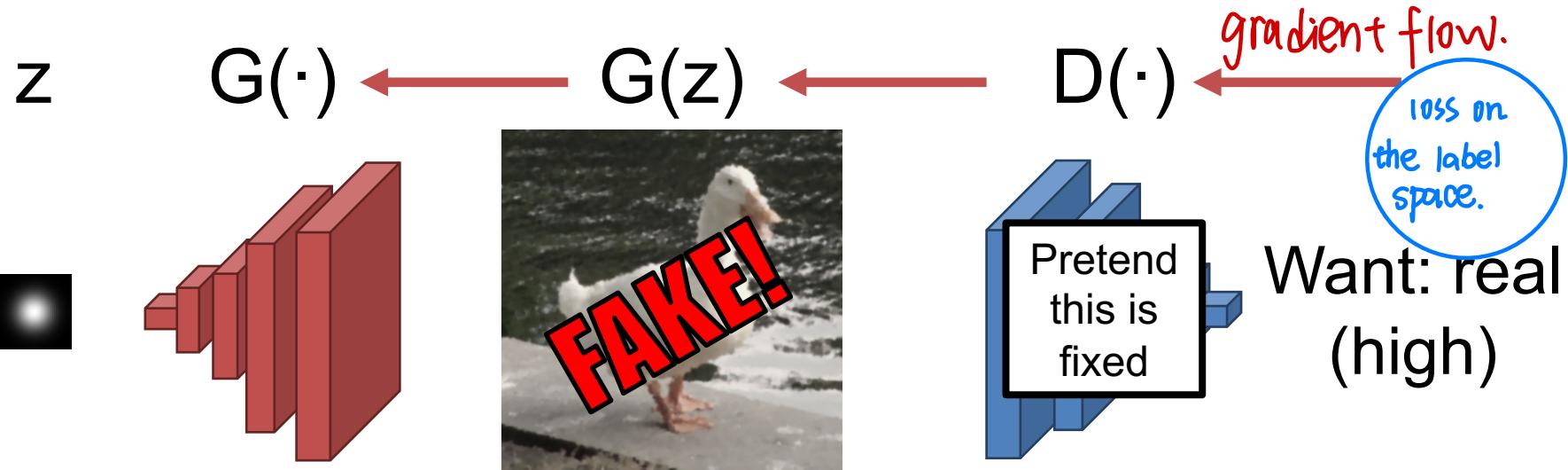
$D(\cdot)$



Want: real  
(high)

$$\arg \max_D E_{z,x} [\log (1 - D(G(z))) + \log(D(x))]$$

# Generative Adversarial Networks



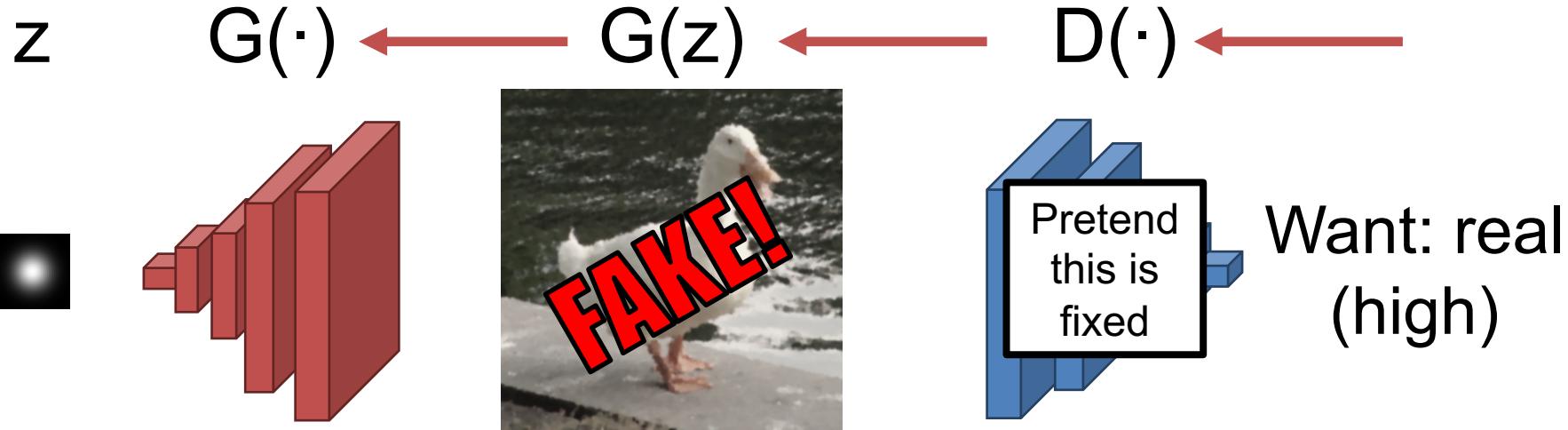
Goal of generator G: fool the discriminator D while getting to use gradients from D

Basically what you do: Analyzing the features learned by D. by classifying this real and fake. and then using these features learned by D to supervise signals to the generator.

$$\arg \min_G E_z [\log (1 - D(G(z)))]$$

$1 - D(G(z)) = p(\text{G}(z) \text{ is fake})$  <sup>minimize</sup>  
Want to maximize  $\log p(\text{fake})$  for fake data

# Generative Adversarial Networks



Goal of generator  $G$ : fool the discriminator  $D$  while getting to use gradients from  $D$

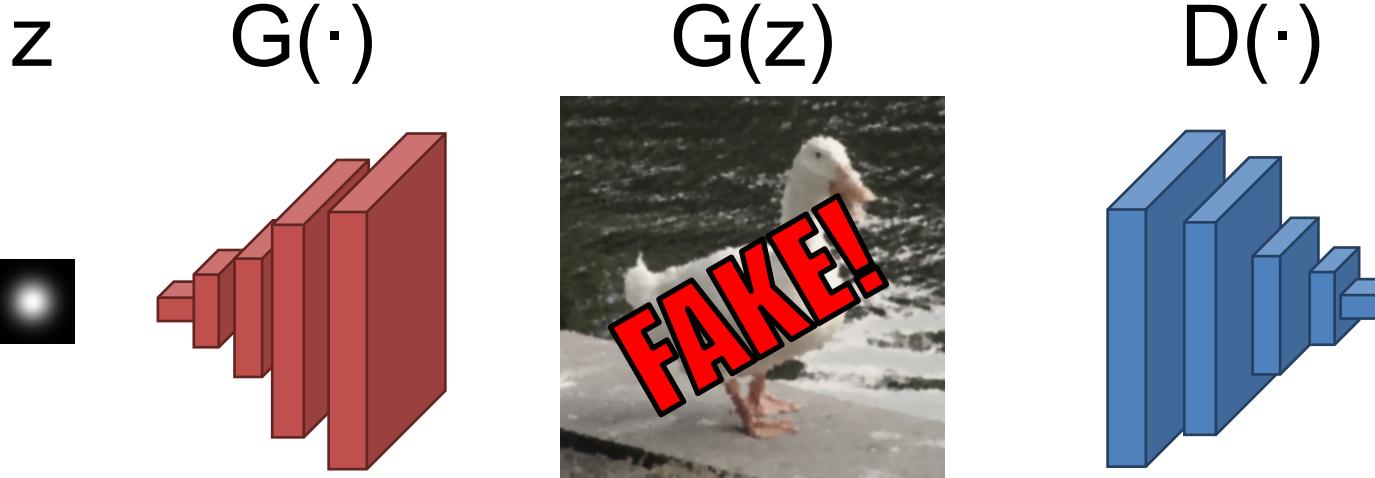
Intuition: make  $D$  think  $G(z)$  is real

Analogy: art forger and art detective

How good are you at spotting forgeries?  
try to fake an image that looks like the original, try to find features that marks the different

$$\arg \min_G E_z [\log (1 - D(G(z)))]$$

# Generative Adversarial Networks



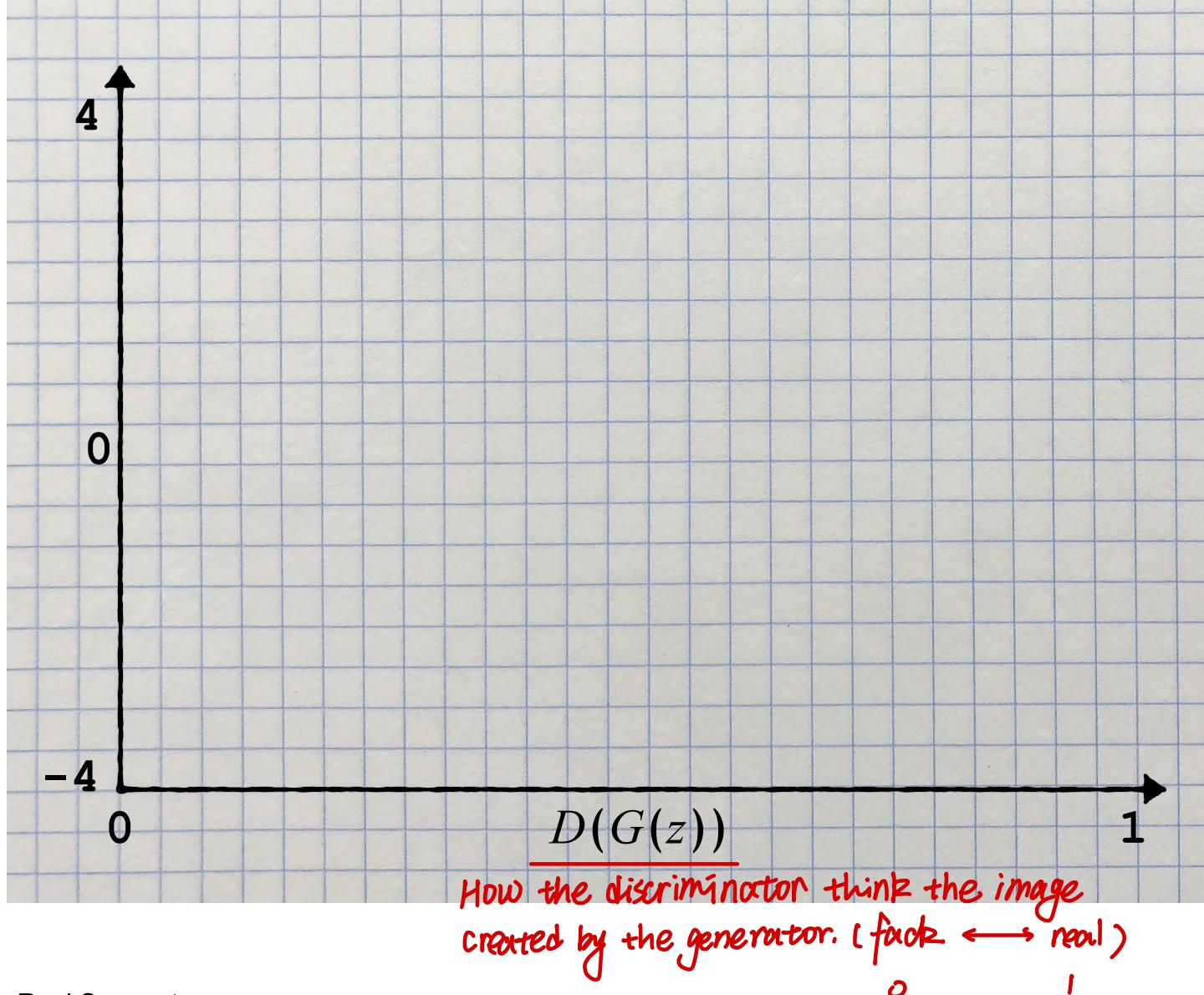
Final goal: find the generator that fools the best D  
that you could find.

Min-max game between G and D

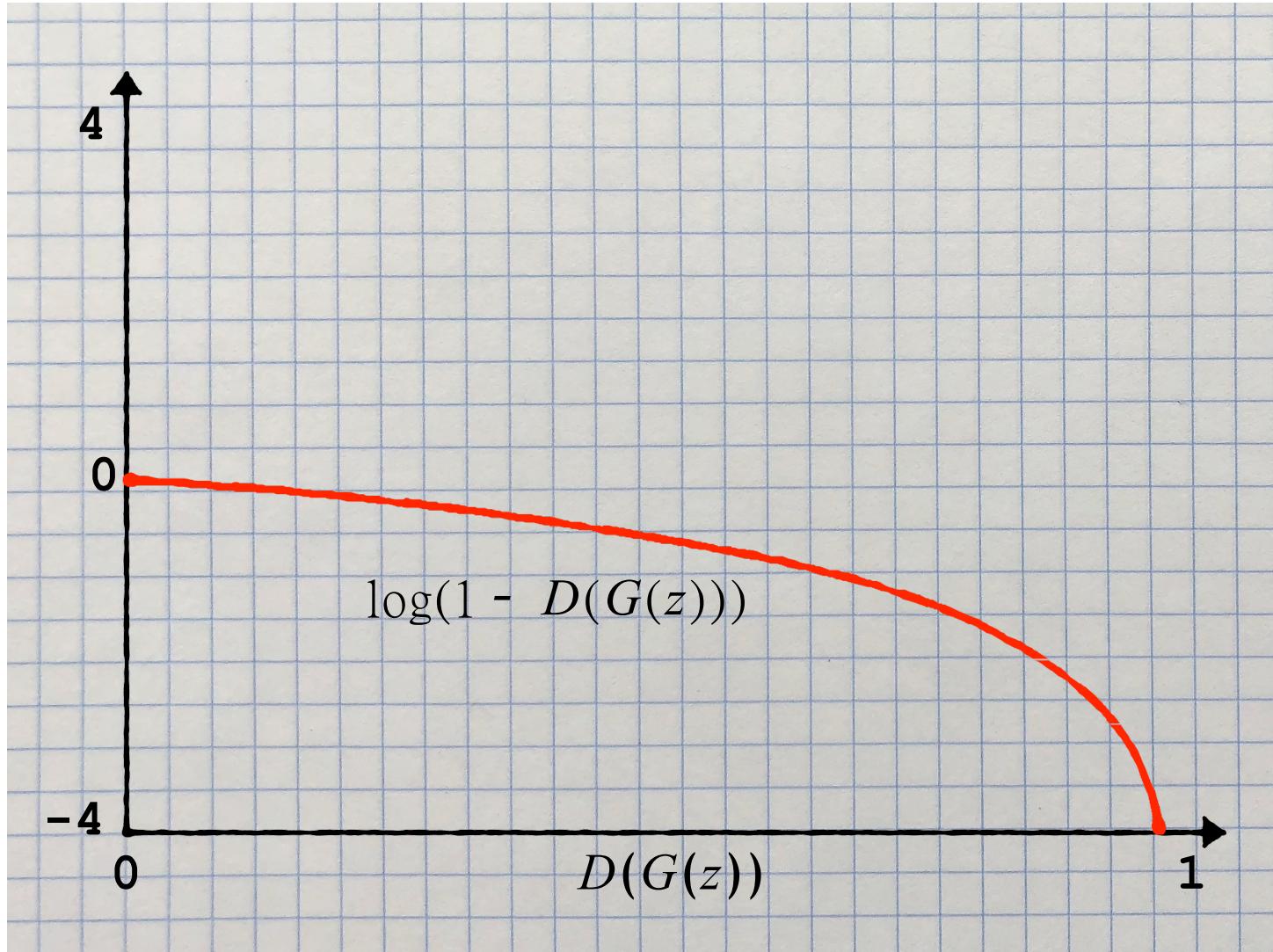
In practice, important not to let the discriminator get  
too good. **Why?**

$$\arg \min_G \max_D E_{z,x} [\log(1 - D(G(z))) + \log(D(x))]$$

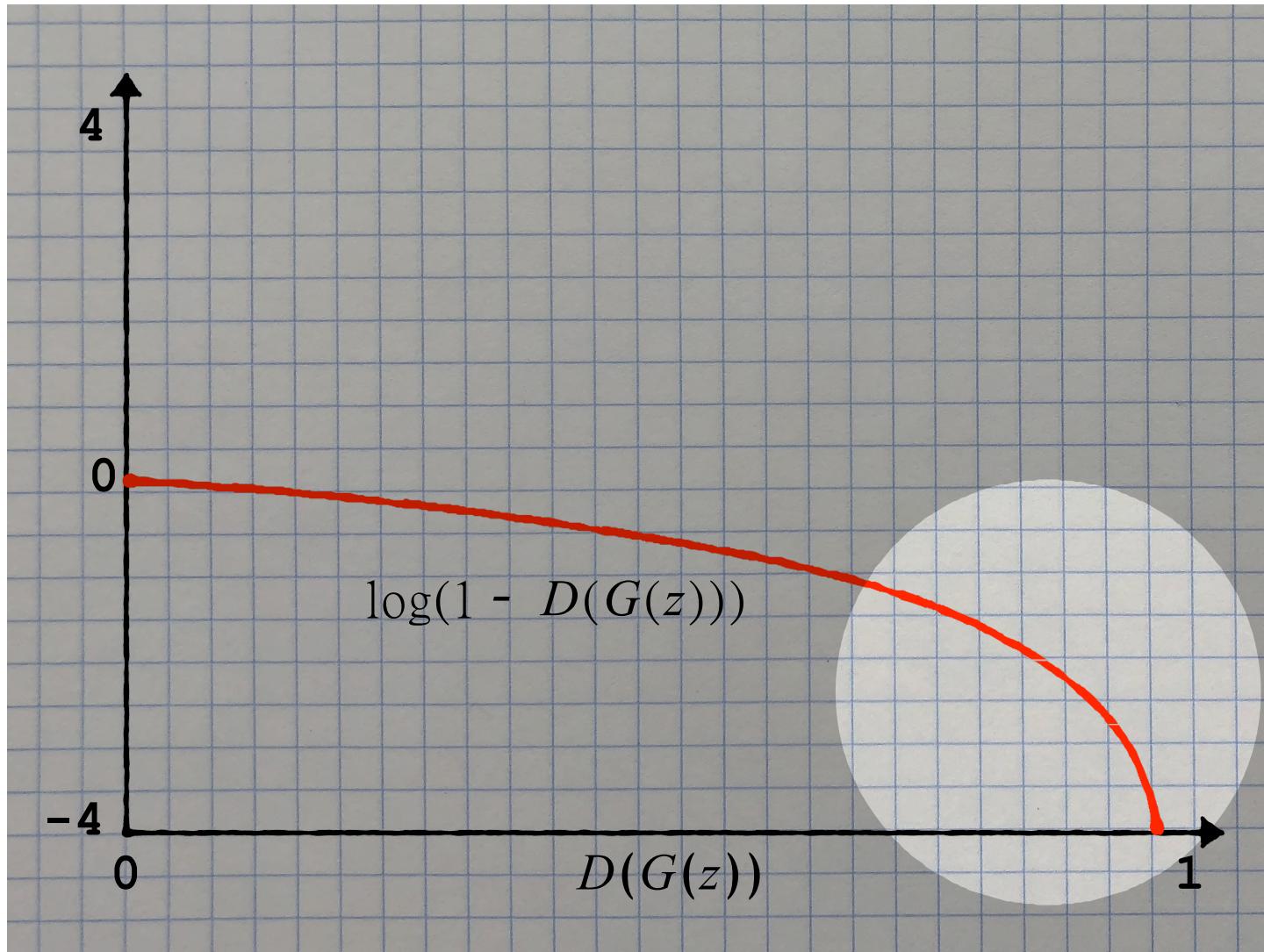
# Caveat on G Loss

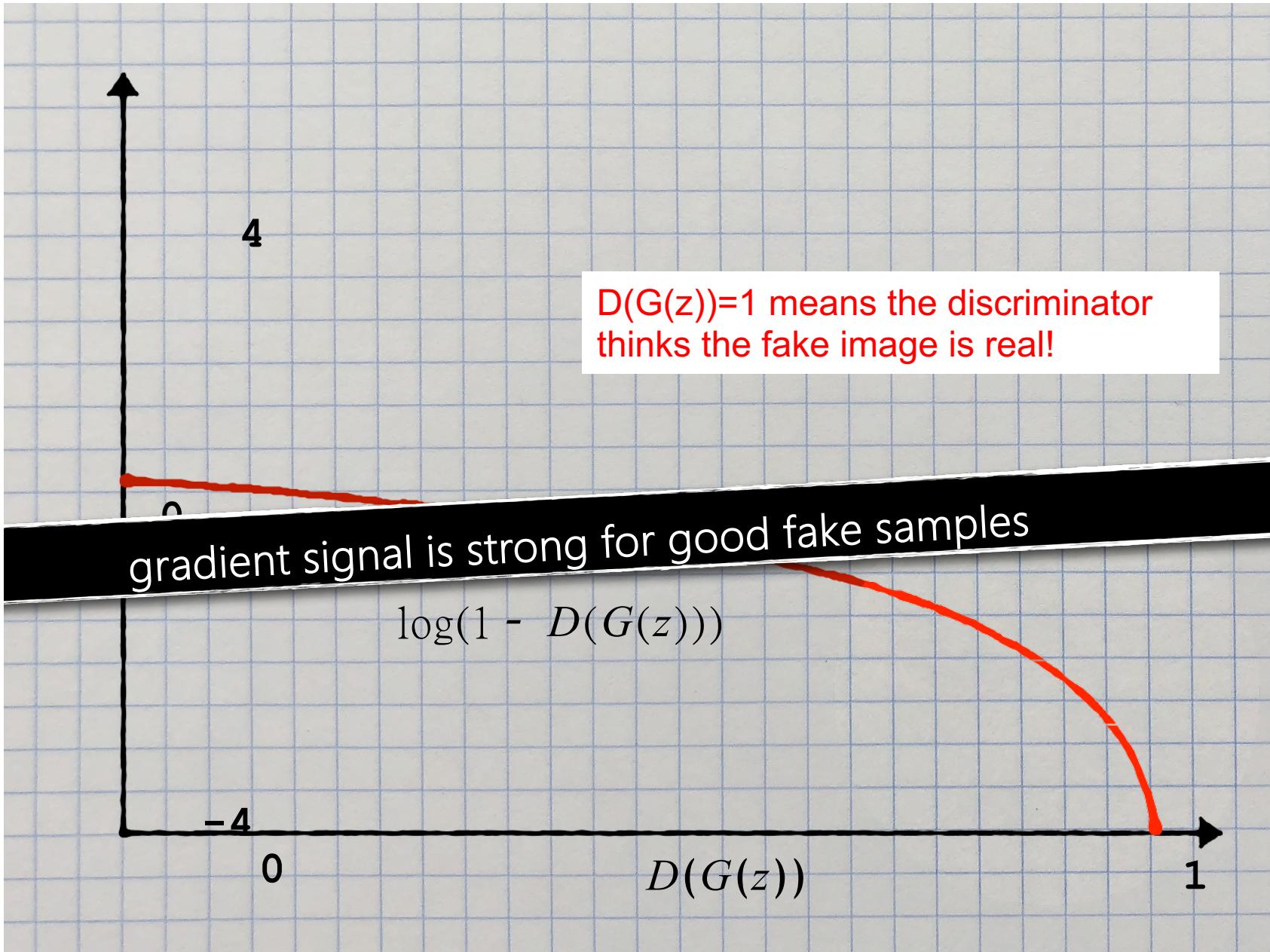


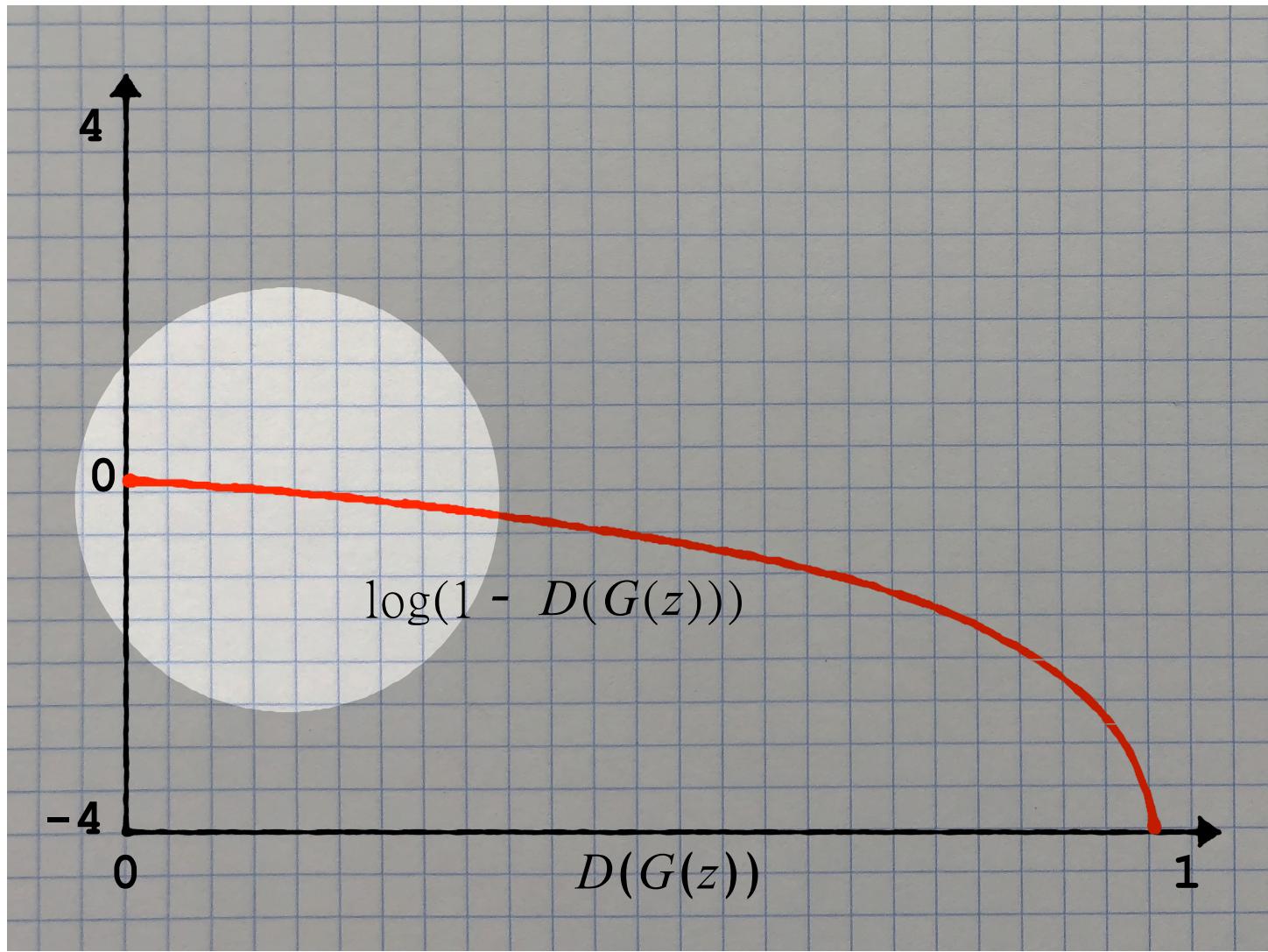
$$\arg \min_G E_z [\log (1 - D(G(z)))]$$



$$\arg \min_G E_z [\log (1 - D(G(z)))]$$







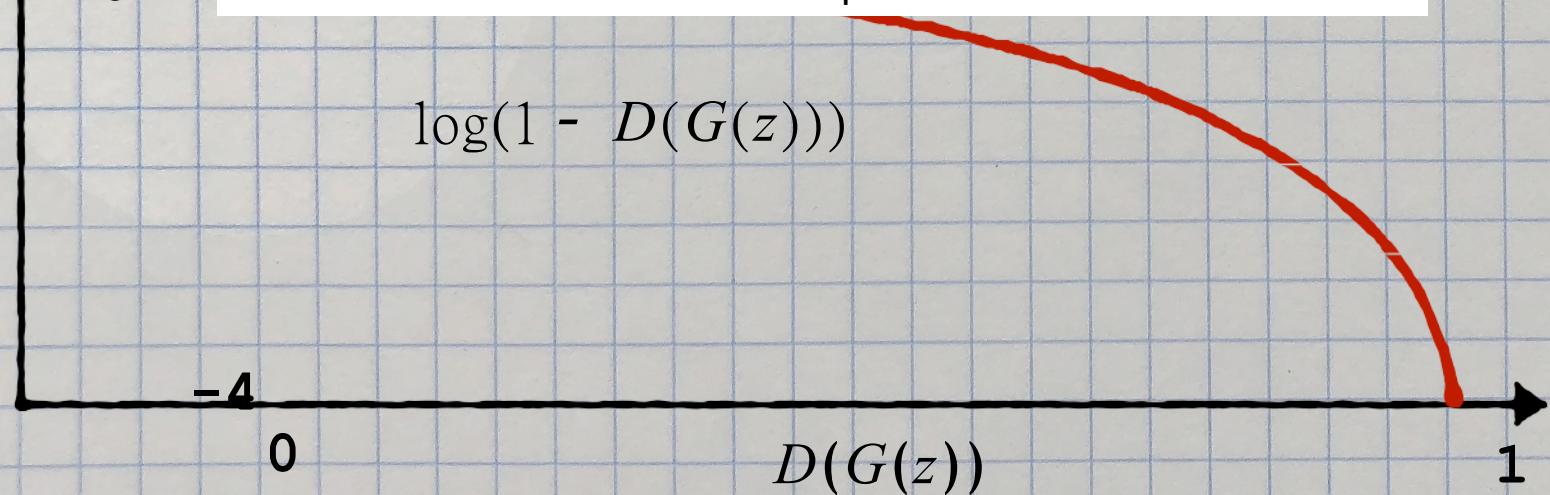


This kind of problem is often known as vanishing gradient

4

gradient signal is weak for bad fake samples

At the beginning of the training most fake samples are bad! Then gradient is small, and the generator do not receive much information from discriminator to update itself!

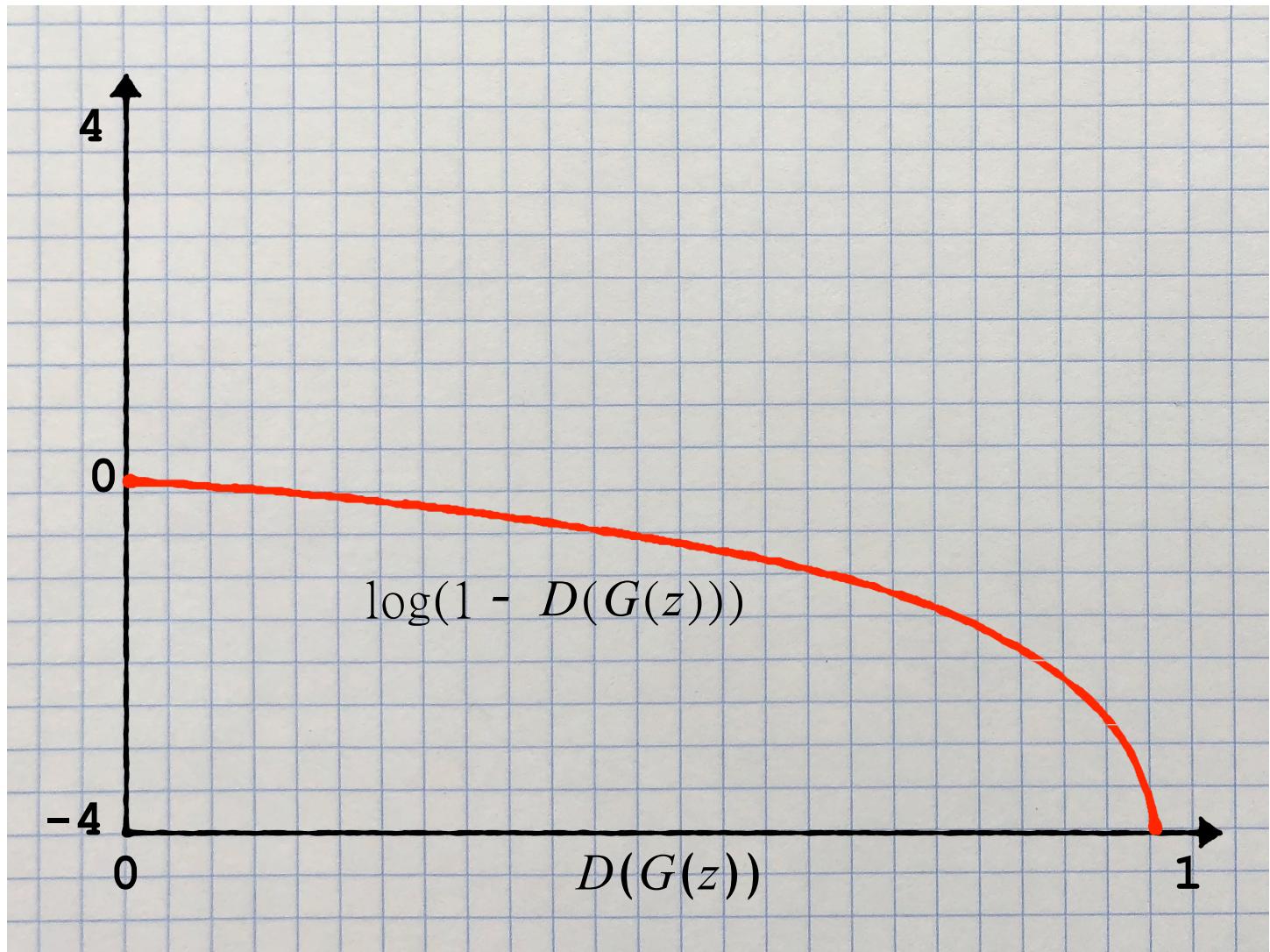


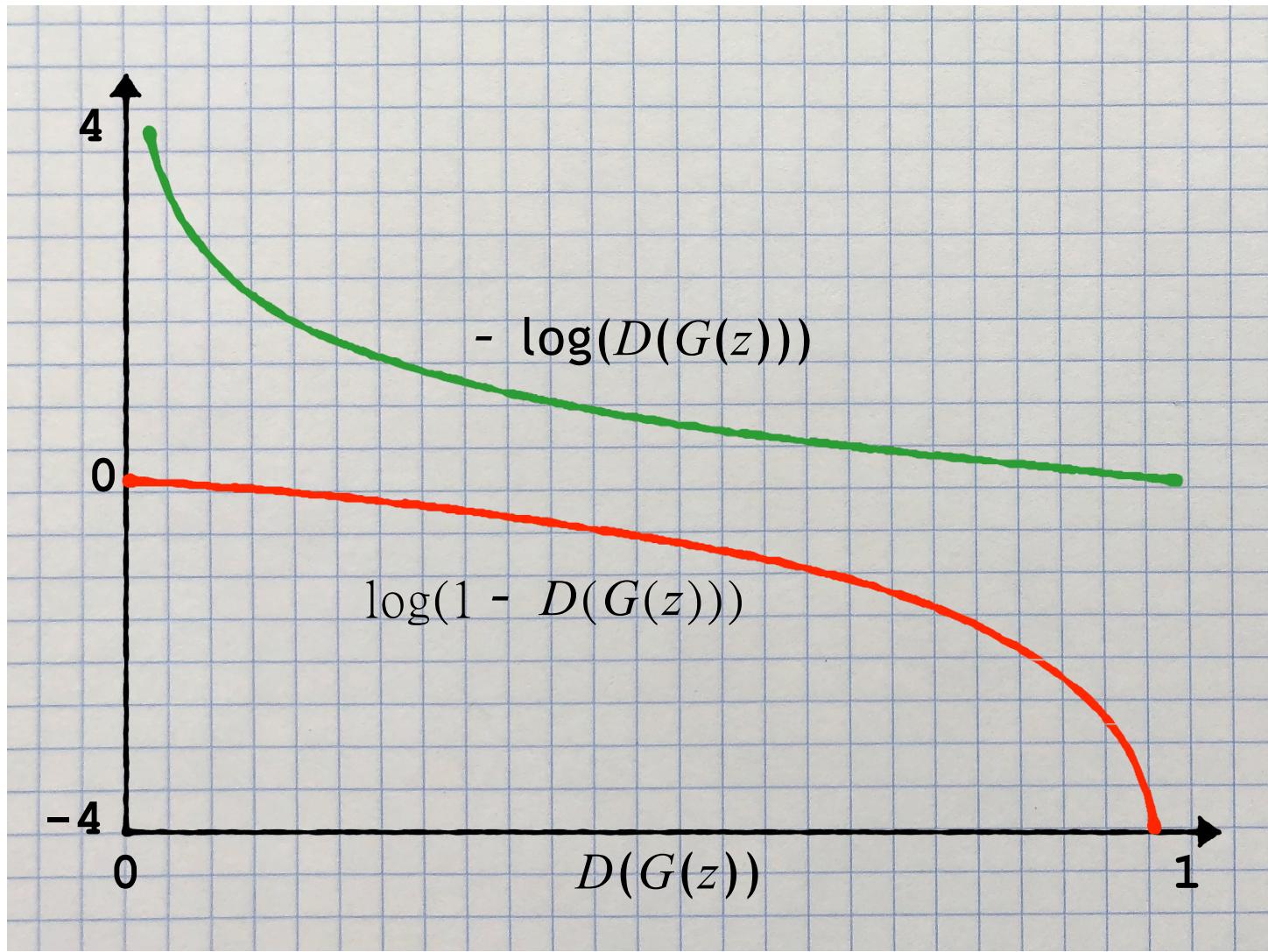
# Improved G Loss

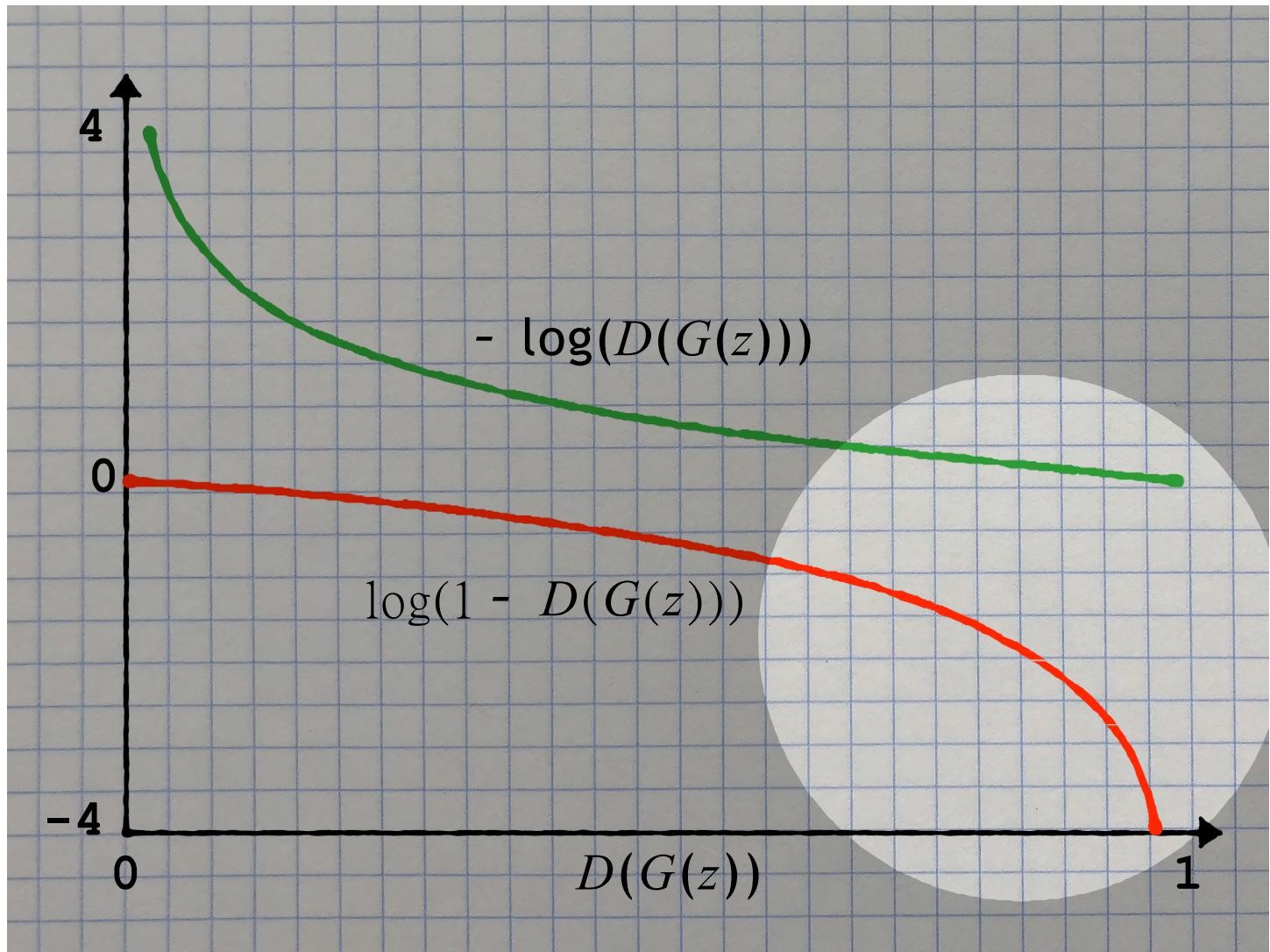
$$E_z[\log(1 - D(G(z)))]$$

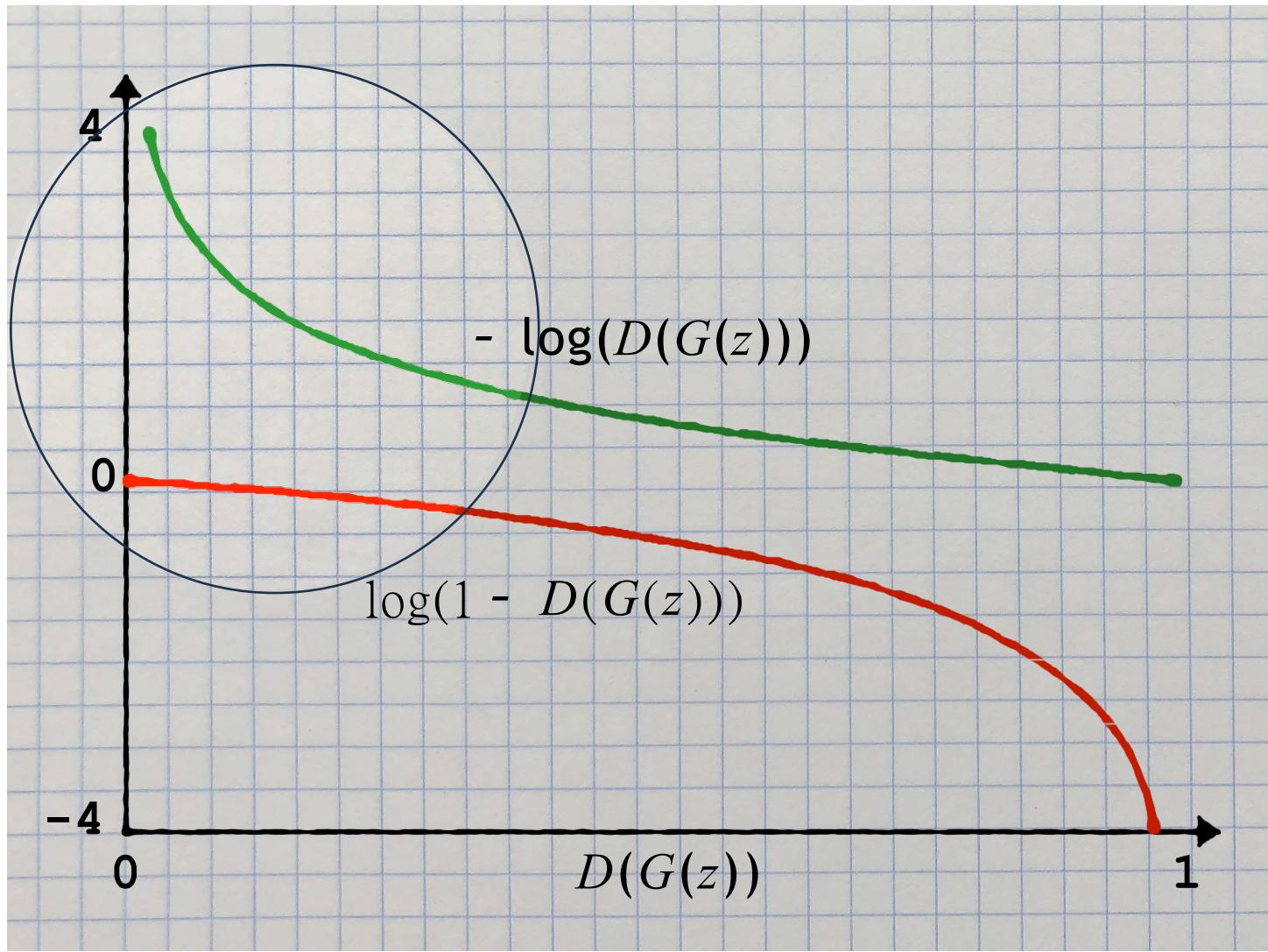
**Replace with**

$$E_z[-\log(D(G(z)))]$$









# Typically for GAN Training

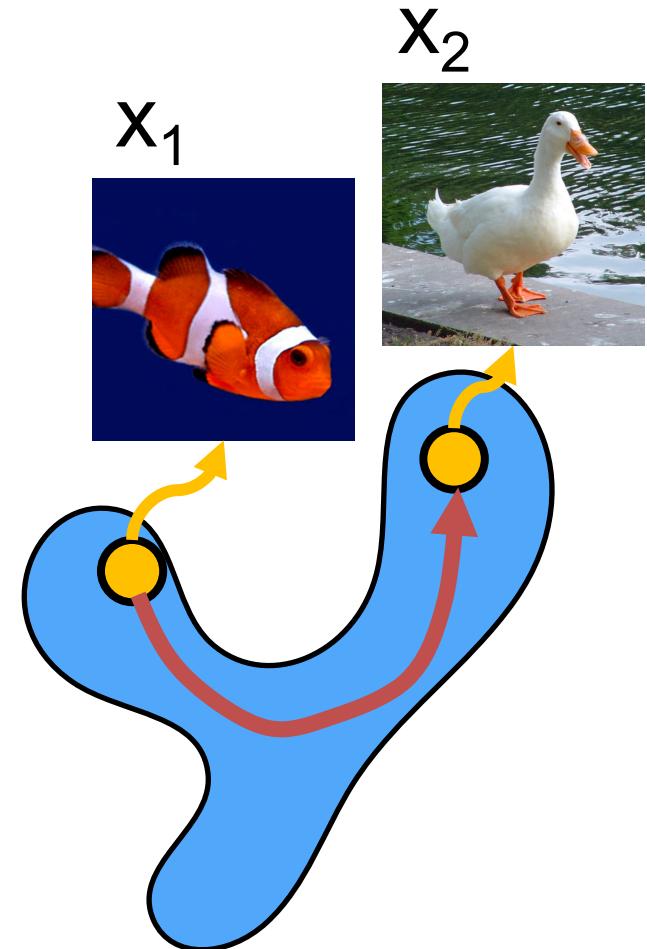
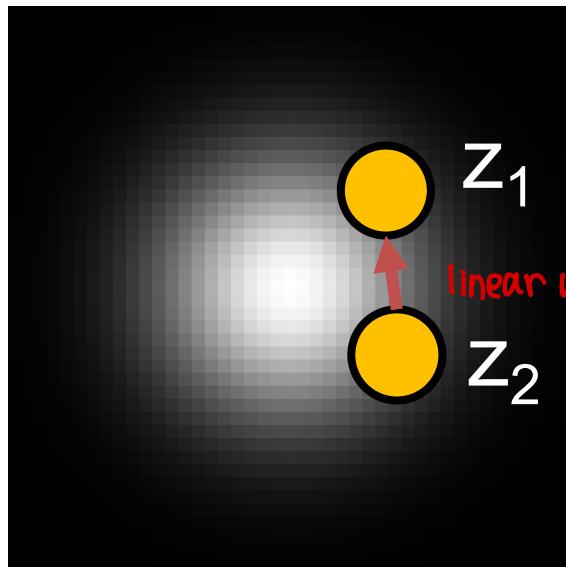
- Very often unstable *since it's complex*
- Need regularizations or more stable losses (out of scope) *to make the discriminator weaker.*
- Need balance between D and G
- Training discriminator is easier than generator (more people can appreciate art vs # of artists)

# GAN Summary

```
# Set batch size
for number of training iterations do
    # D update
    Sample noise vectors z from noise prior
    Generate images G(z)
    Sample real images x
    Compute D loss:  $\ell_D = \sum -\log(1 - D(G(z))) - \log(D(x))$ 
    Take gradient descent step to update D
    # G update
    Sample noise vectors z from noise prior
    Generate images G(z)
    Compute G loss:  $\ell_D = \sum -\log(D(G(z)))$ 
    Take gradient descent step to update G
end for
```

# Revisiting Averages

Can use  $z$  to walk latent space  
 $G(\alpha z_1 + (1 - \alpha) z_2)$  for  $\alpha$  in  $[0, 1]$



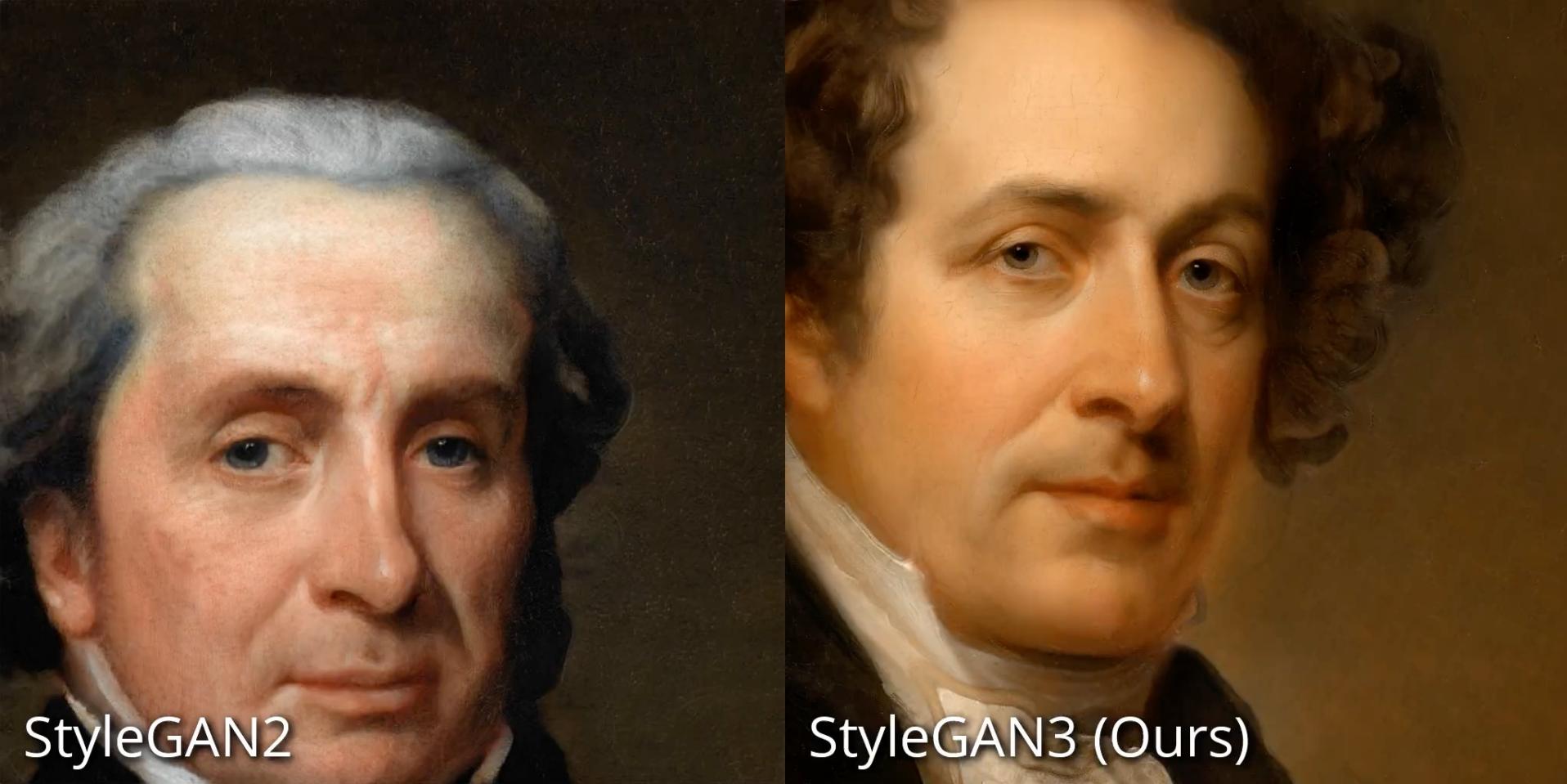


StyleGAN2



StyleGAN3 (Ours)

[Karras et al., “Alias-Free Generative Adversarial Networks”, 2021]



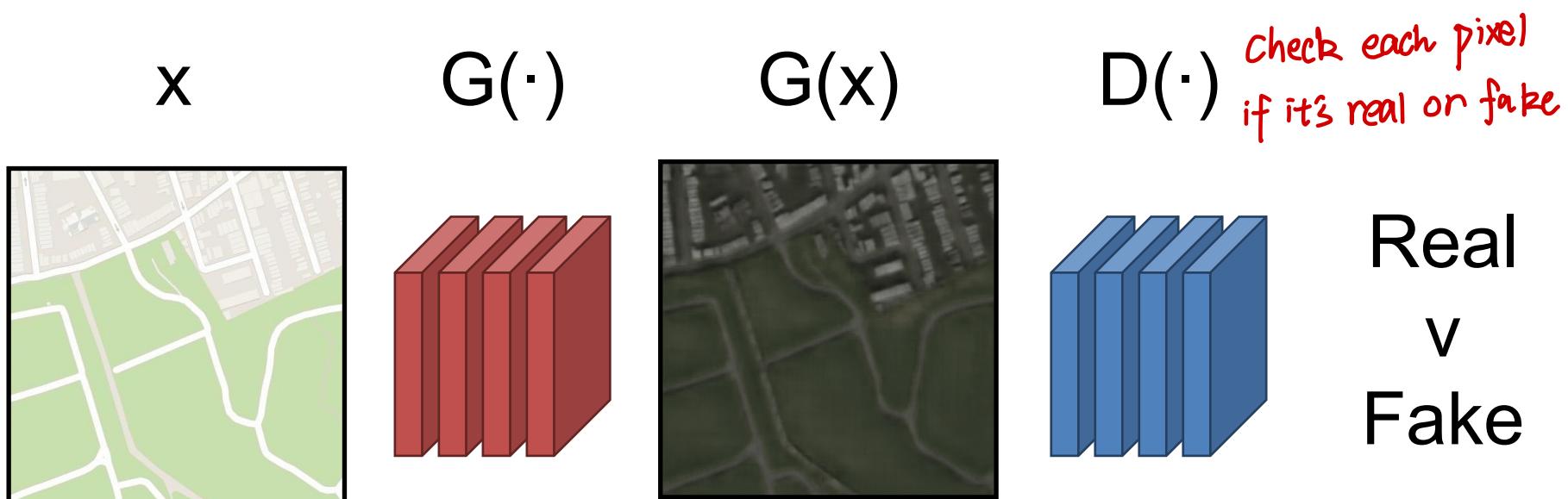
[Karras et al., “Alias-Free Generative Adversarial Networks”, 2021]

# Conditioning on Things

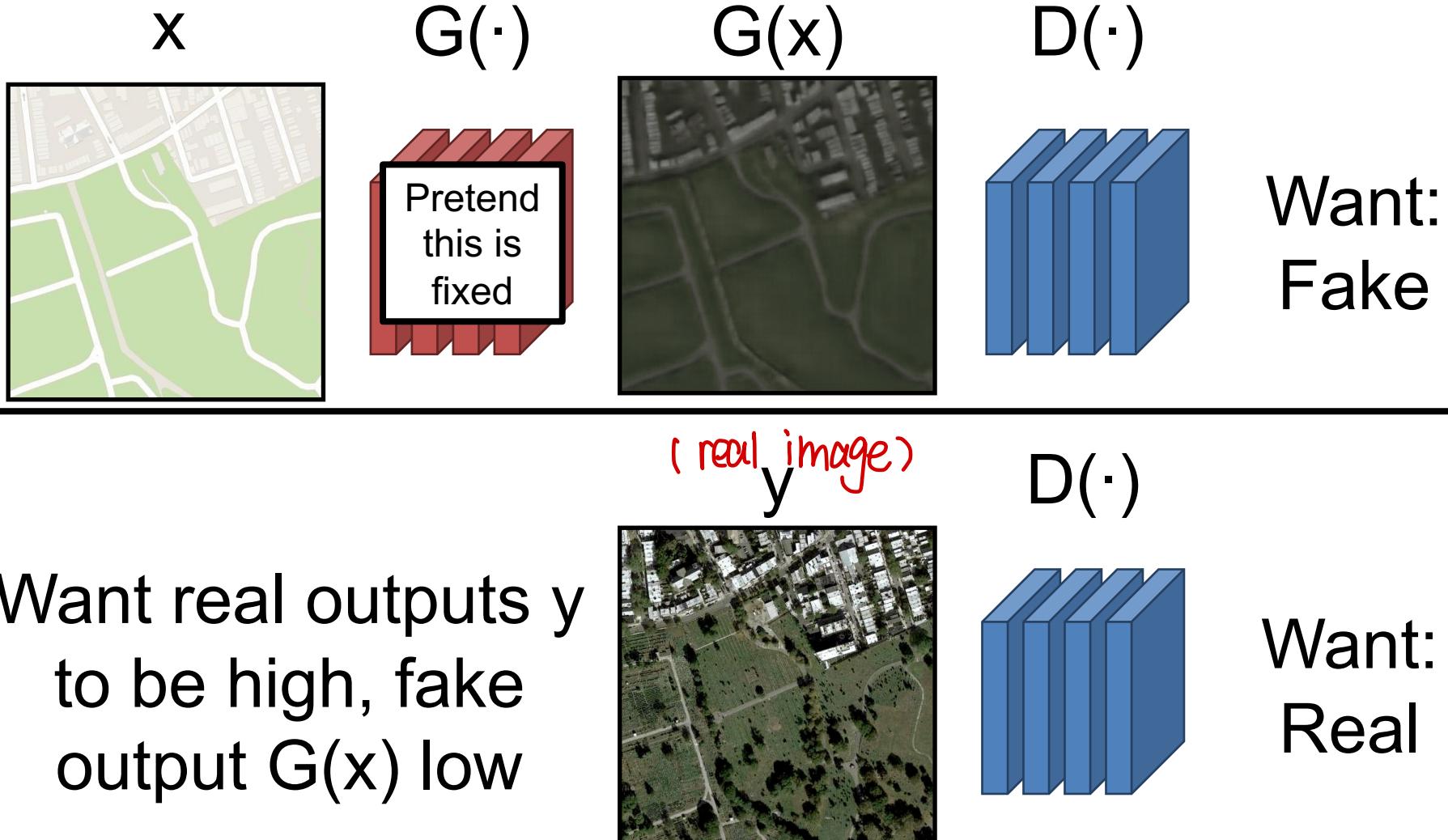
- Turning noise into pictures of things is all fun and good, but what if we want control over our synthetic images?

# Conditional GANs (Pix2Pix)

- Generator tries to make fake images – accepts **image** and makes an image
- Discriminator tries to identify fakes – outputs  $p(\text{fake})$ , potentially at each pixel

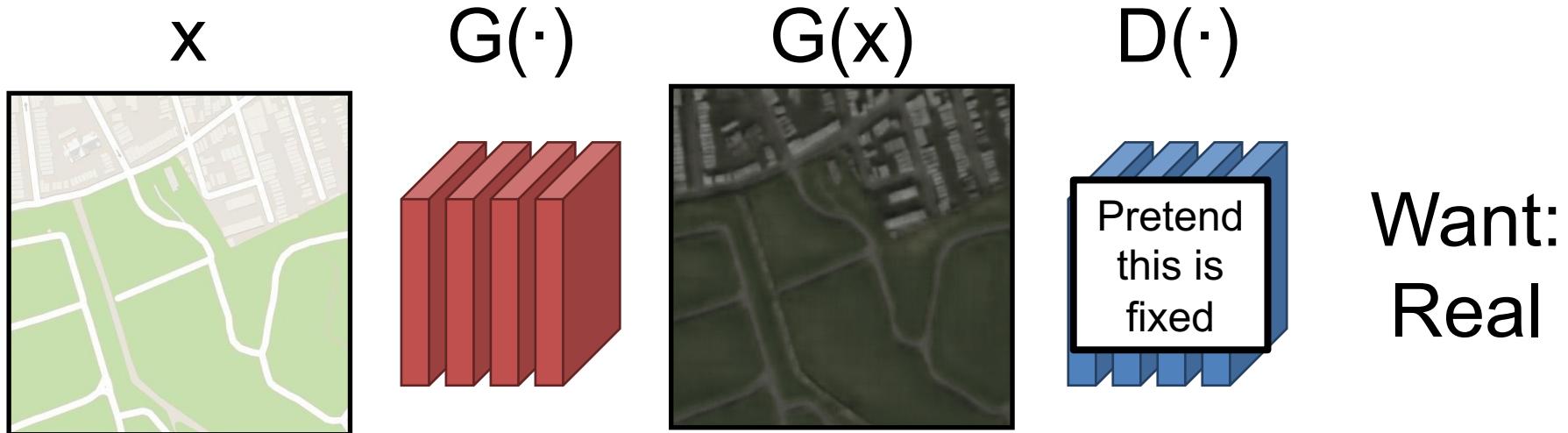


# Conditional GANs – Discriminator



$$\arg \max_D E_{z,x} [ \log(1 - D(G(x))) + \log(D(y)) ]$$

# Conditional GANs – Generator



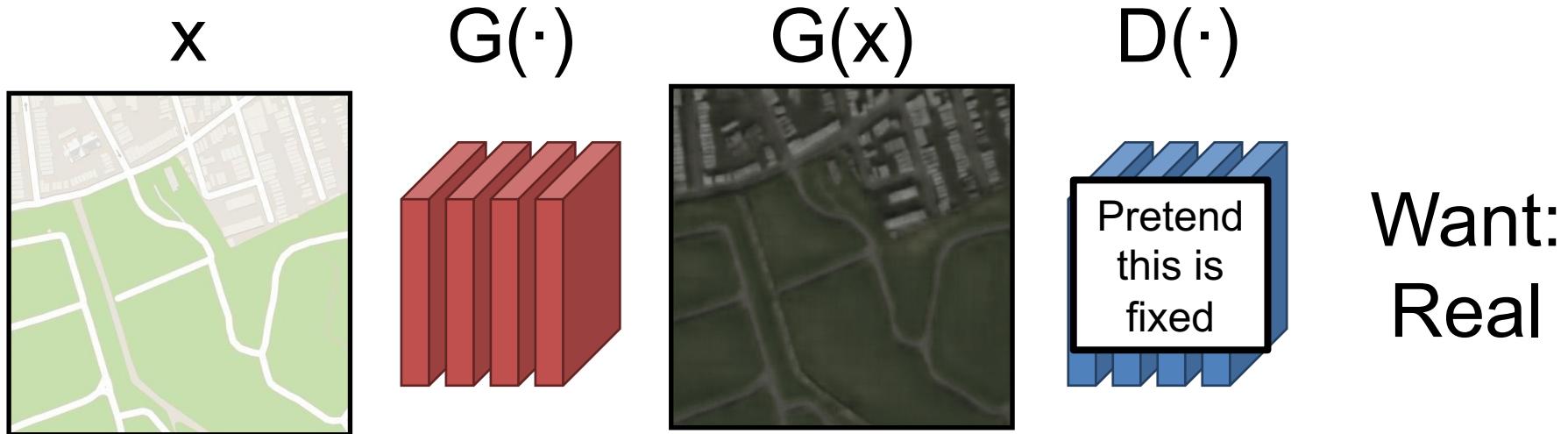
If you're the generator, want to make fakes that fool the discriminator into think they're real

$$\arg \min_G E_{z,x} [\log(1 - D(G(x)))]$$

Same min/max game as before

$$\arg \min_G \max_D E_{z,x} [\log(1 - D(G(x))) + \log(D(y))]$$

# Conditional GANs – Generator



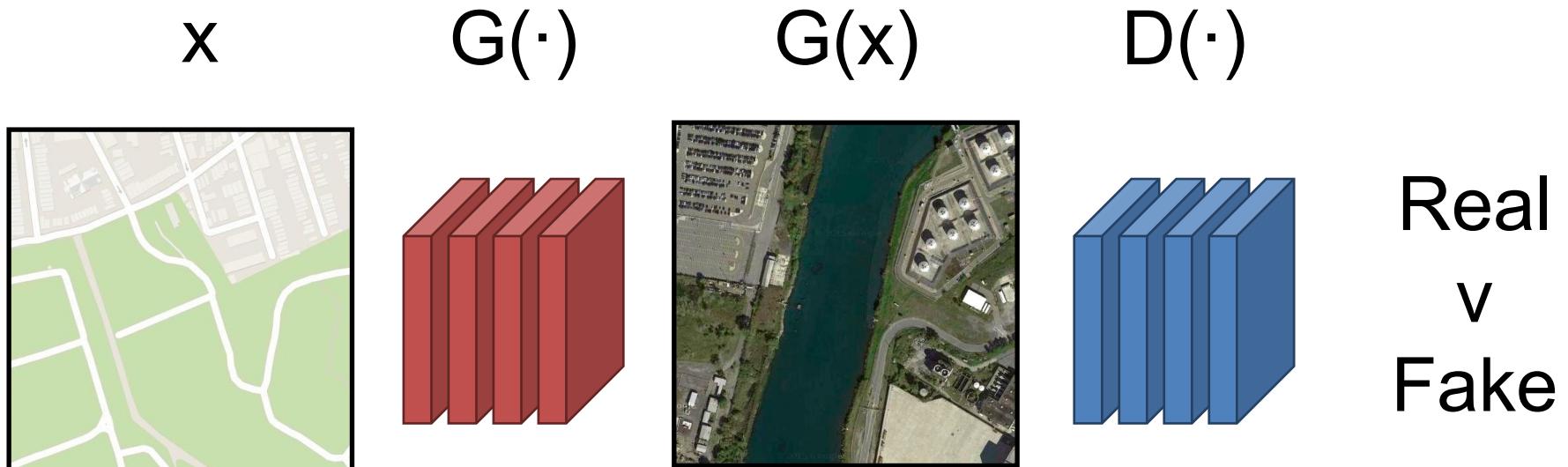
In practice, use alternating optimization  
using the improved G loss

$$\arg \min_G E_{z,x}[-\log(D(G(x)))]$$

$$\arg \min_D E_{z,x}[-\log(1 - D(G(x))) - \log(D(y))]$$

# One Catch

- $G$  can just output random good images.
- Solution – additional L1 loss

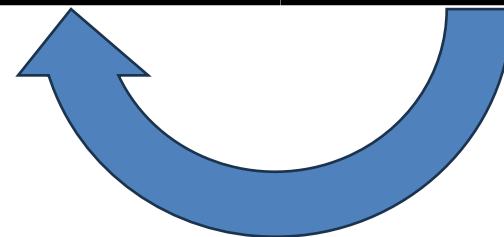
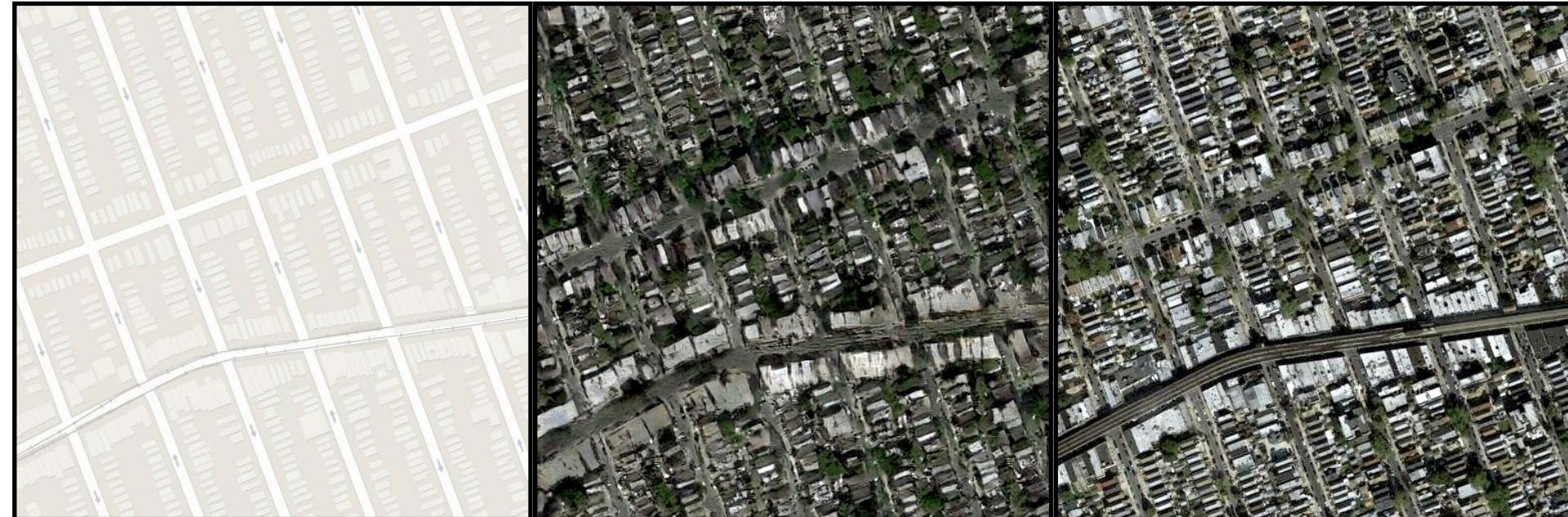


# GAN + L1 Loss

Input

Output

Groundtruth

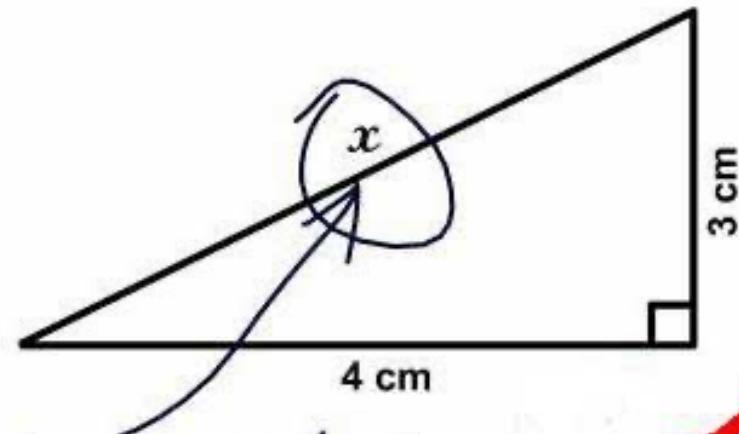


Pixelwise L1 Loss

# More Broadly

- Neural networks are lazy and will do precisely what you ask and no more
- You *have to* be careful what you ask them to do

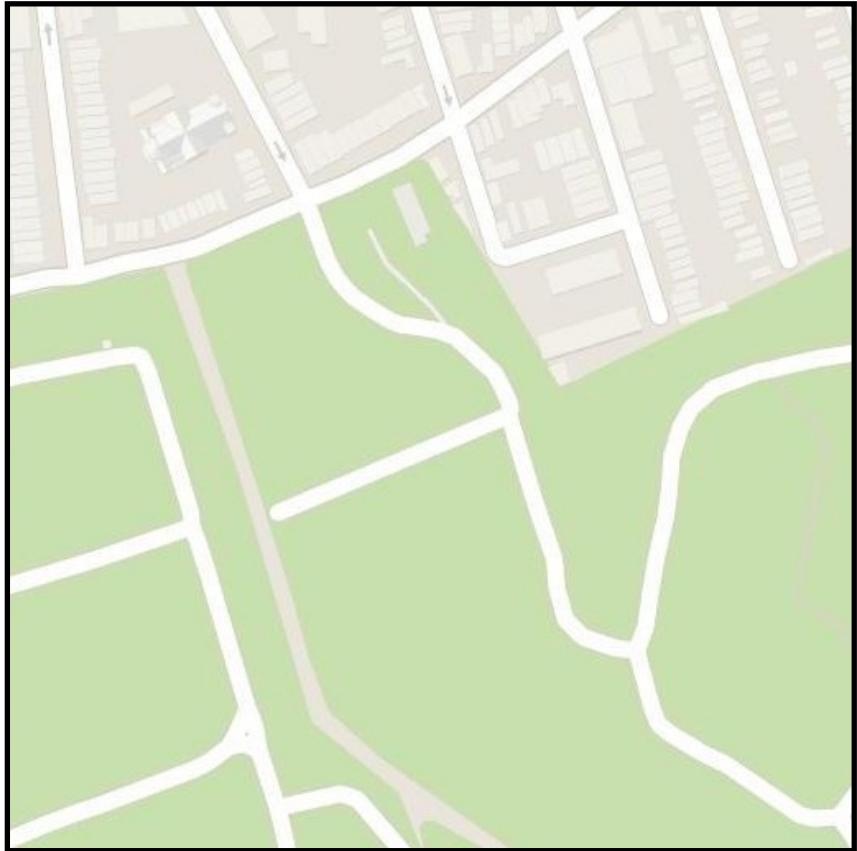
3. Find  $x$ .



Here it is

✓  
-1

Input

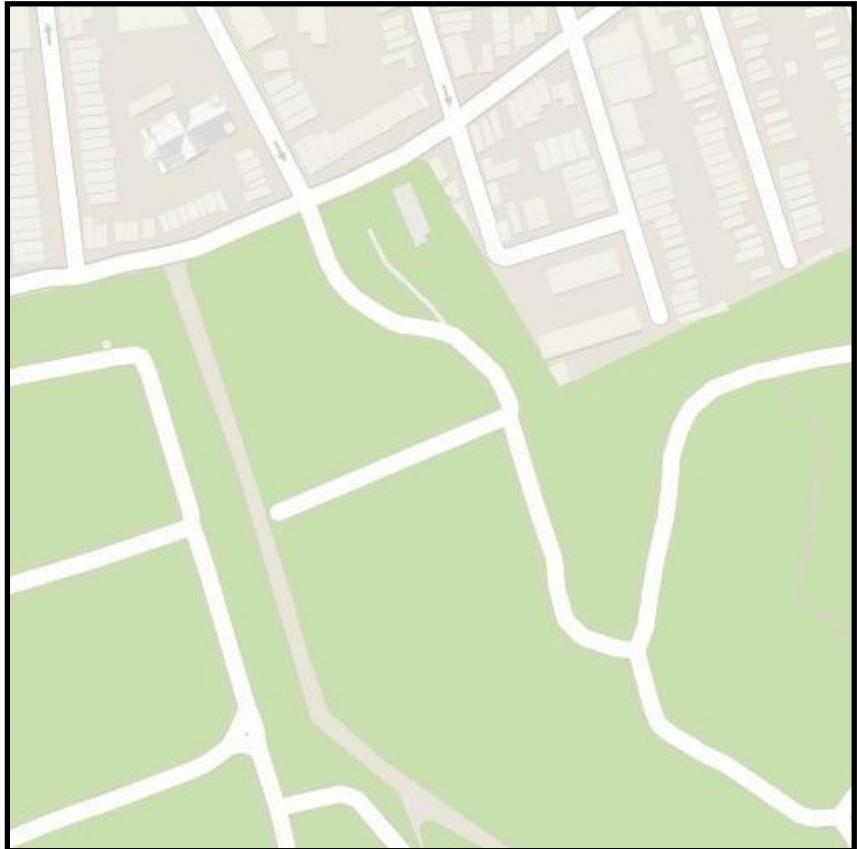


L1 loss



Why is it blurry junk? Hold that thought!

Input



L1 loss + discriminator  
*gan loss*



# Image-to-Image Translation with Conditional Adversarial Networks

Phillip Isola

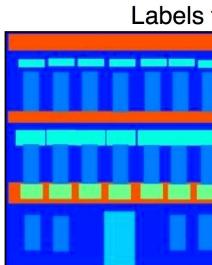
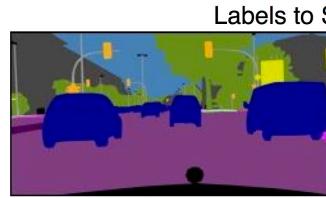
Jun-Yan Zhu

Tinghui Zhou

Alexei A. Efros

Berkeley AI Research (BAIR) Laboratory, UC Berkeley

{isola, junyanz, tinghuiz, efros}@eecs.berkeley.edu



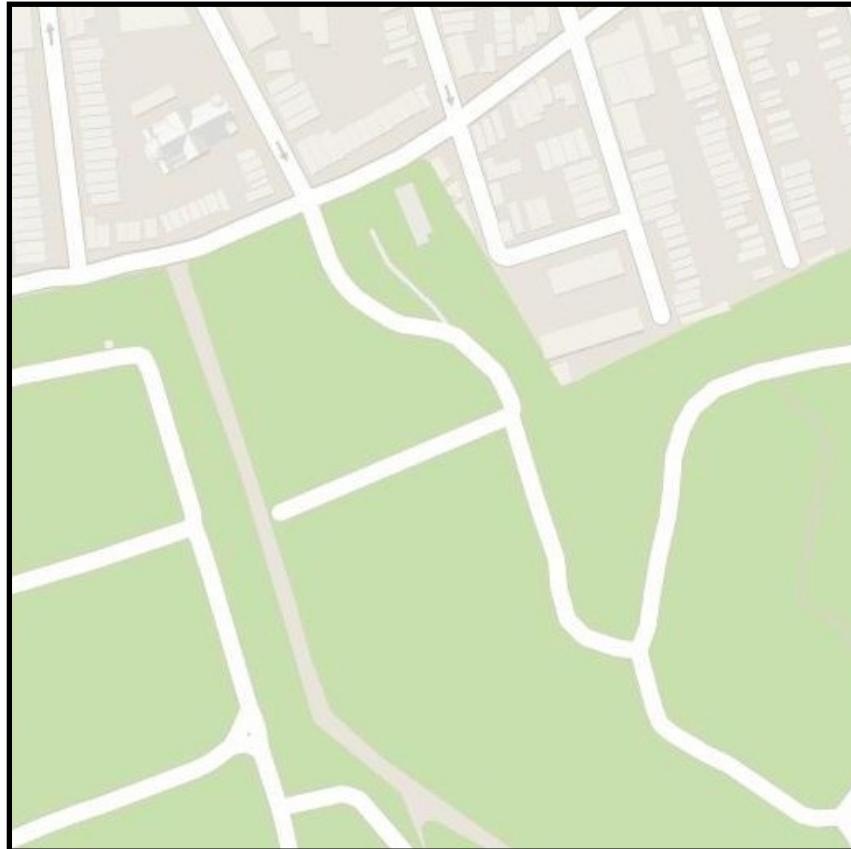
BW to Color



Edges to Photo



# Let's Talk About Blurry Pictures

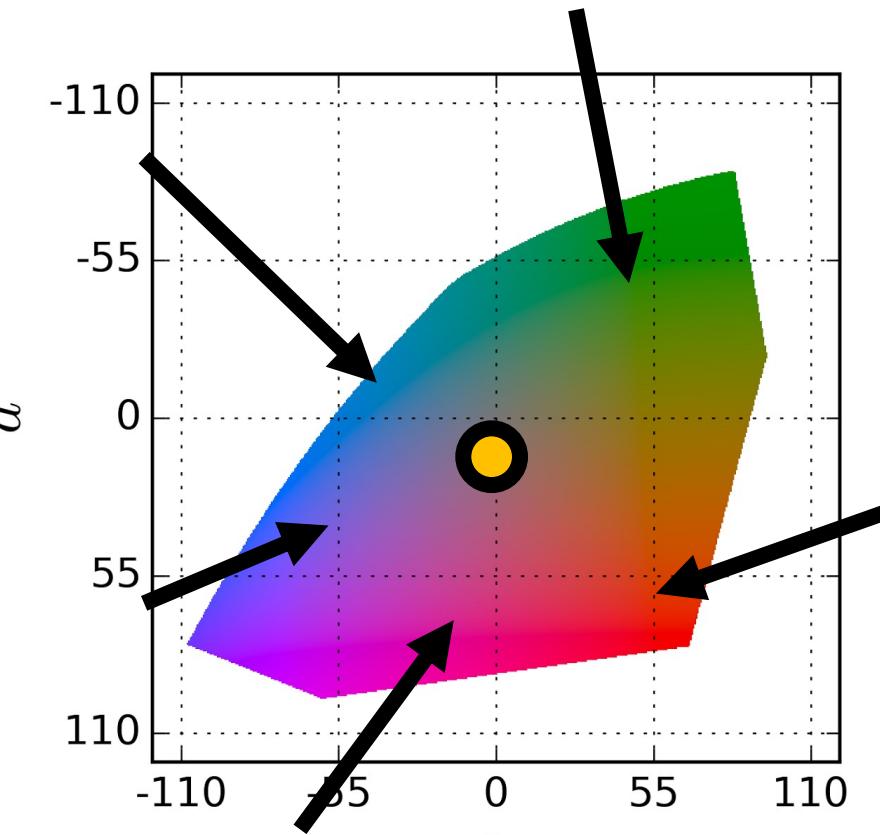


# What Color Is This Bird?



To make things more concrete: what color is this the pixel under this gold circle?

# What Color Is This Bird?



$L_2$  : average

$L_1$  : median.

Many options. What minimizes mean-squared error?  
What minimizes the L1 distance?

# What Color Is This Bird?

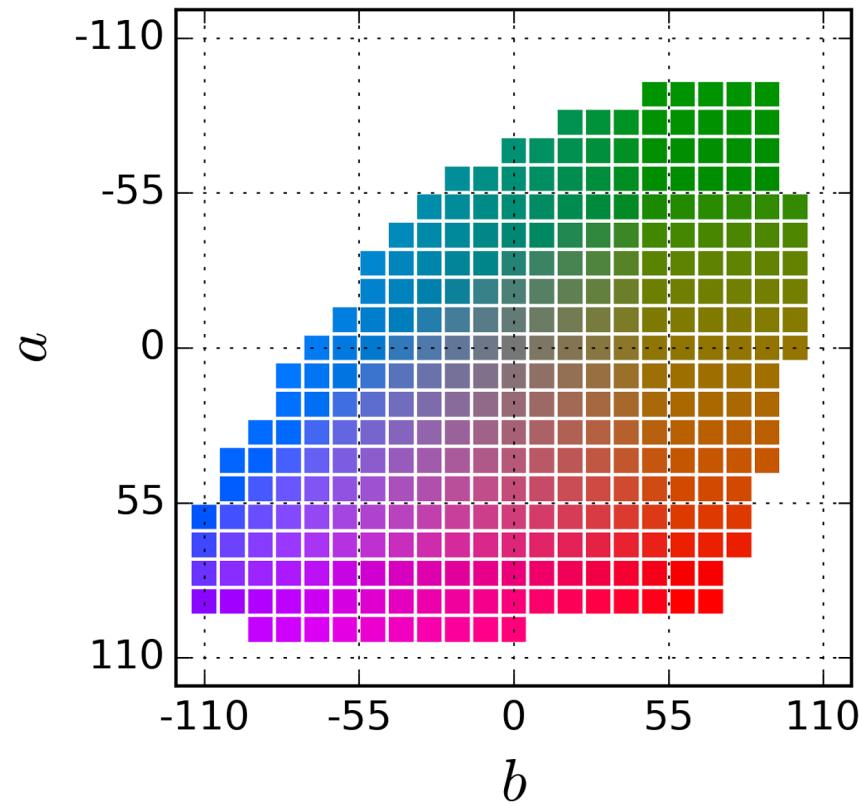
Option 1: Discretize / quantize

*Before learning:* assign pixel nearest color index

*After:* convert index to value

Works because network can express its uncertainty.

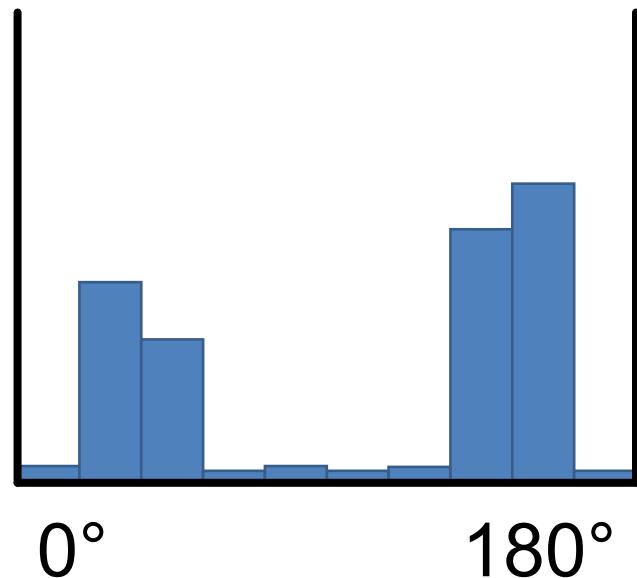
Intuition: RGB prediction as a classification task



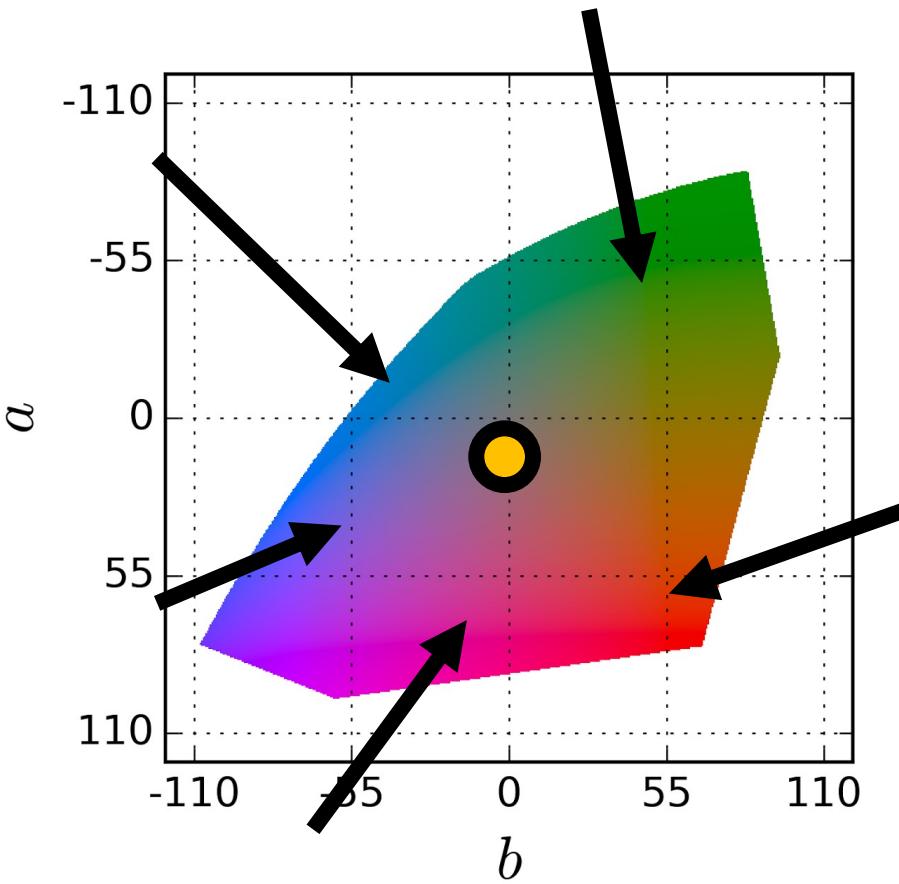
# Discretized Values – Angles

Imagine predicting an angle from  $0^\circ$  to  $180^\circ$ . Having bins enables:

- Expressing bimodal distributions (e.g., either  $30^\circ$  or  $150^\circ$ )
- Getting a confidence from the prediction



# Option 2 – GAN



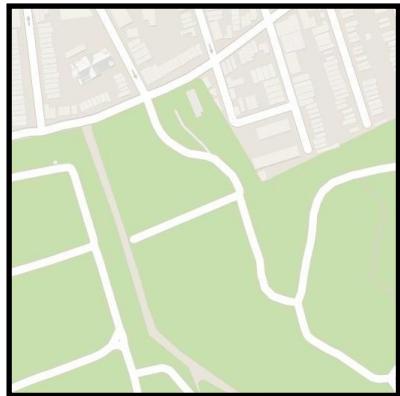
GAN – the discriminator will prevent us from making birds grey or brown. **Why?**

# What to Take Away

- Be careful what you ask a deep net to solve.
- The objective you're asking it to solve bakes in assumptions
- Most solutions broken in one way or another
- Deep learning is not magic

# Aside: Perceptual Losses

$G(\cdot)$



$G(x)$



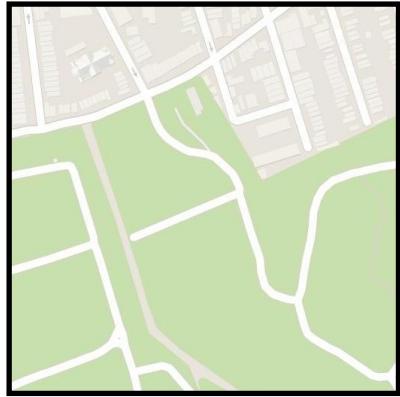
$y$

Conventionally,  
minimize distance in  
pixel space:  
 $\|G(x) - y\|$



# Aside: Perceptual Losses

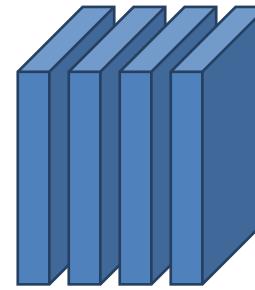
$G(\cdot)$



$G(x)$



$F(\cdot)$



$F(G(x))$



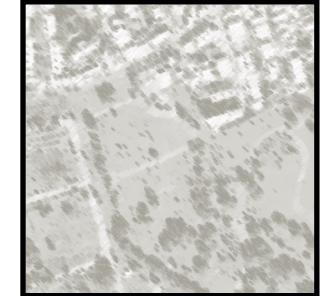
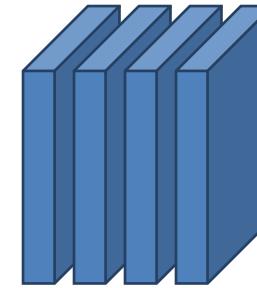
$y$

Instead measure  
distance after passing  
through pre-trained  
network

$$\|F(G(x)) - F(y)\|$$



$F(y)$



# CycleGAN

- Conditional GAN but don't have paired data
- Only have two unpaired datasets



# CycleGAN

Monet  $\curvearrowright$  Photos



Monet  $\rightarrow$  photo

Zebras  $\curvearrowright$  Horses



zebra  $\rightarrow$  horse

Summer  $\curvearrowright$  Winter



summer  $\rightarrow$  winter

photo  $\rightarrow$  Monet



horse  $\rightarrow$  zebra



winter  $\rightarrow$  summer

Photograph



Monet



Van Gogh



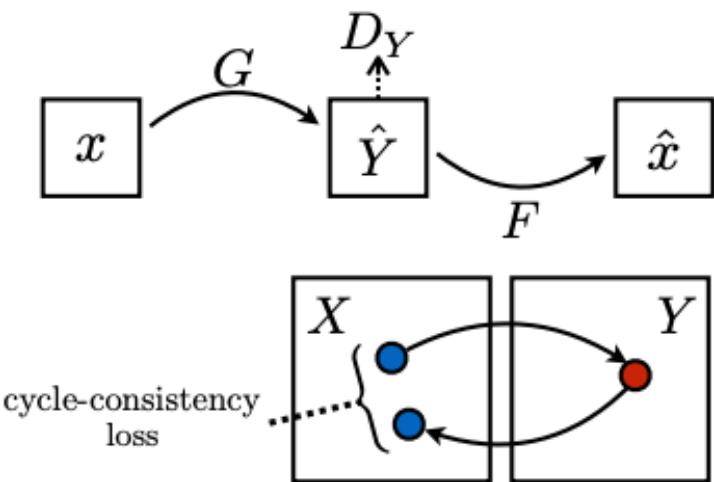
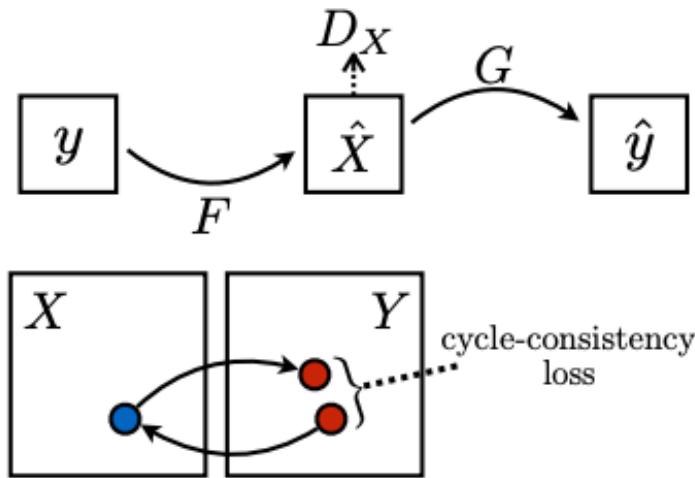
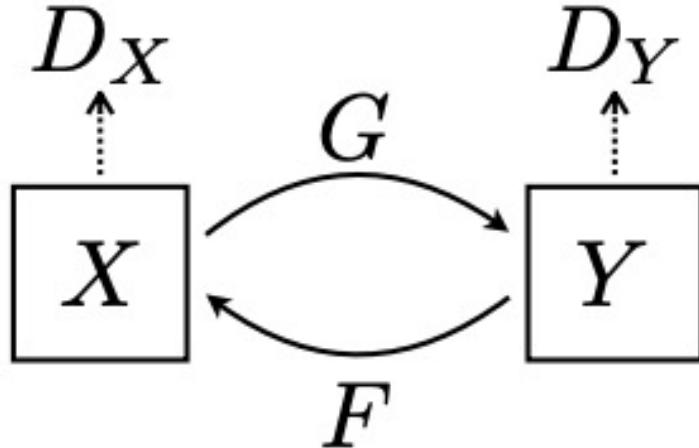
Cezanne



Ukiyo-e

Training data: A set of images of style X + A set of images of style Y

Test: Given an image of style X, generate the same image in style Y



# Unpaired Image Translation

Conditional generative model  $P(\text{ zebra images} | \text{ horse images})$



Zhu et al., 2017

# Next Class

- Diffusion Models